



Titre: Modelling Environmental Effect Dependencies with Principal
Title: Component Analysis and Bayesian Dynamic Linear Models

Auteur: Catherine Paquin
Author:

Date: 2018

Type: Mémoire ou thèse / Dissertation or Thesis

Référence: Paquin, C. (2018). Modelling Environmental Effect Dependencies with Principal
Citation: Component Analysis and Bayesian Dynamic Linear Models [Master's thesis, École Polytechnique de Montréal]. PolyPublie. <https://publications.polymtl.ca/3294/>

 **Document en libre accès dans PolyPublie**
Open Access document in PolyPublie

URL de PolyPublie: <https://publications.polymtl.ca/3294/>
PolyPublie URL:

**Directeurs de
recherche:** James Goulet
Advisors:

Programme: Génie civil
Program:

UNIVERSITÉ DE MONTRÉAL

MODELLING ENVIRONMENTAL EFFECT DEPENDENCIES WITH PRINCIPAL
COMPONENT ANALYSIS AND BAYESIAN DYNAMIC LINEAR MODELS

CATHERINE PAQUIN
DÉPARTEMENT DES GÉNIES CIVIL, GÉOLOGIQUE ET DES MINES
ÉCOLE POLYTECHNIQUE DE MONTRÉAL

MÉMOIRE PRÉSENTÉ EN VUE DE L'OBTENTION
DU DIPLÔME DE MAÎTRISE ÈS SCIENCES APPLIQUÉES
(GÉNIE CIVIL)
AOÛT 2018

UNIVERSITÉ DE MONTRÉAL

ÉCOLE POLYTECHNIQUE DE MONTRÉAL

Ce mémoire intitulé :

MODELLING ENVIRONMENTAL EFFECT DEPENDENCIES WITH PRINCIPAL
COMPONENT ANALYSIS AND BAYESIAN DYNAMIC LINEAR MODELS

présenté par : PAQUIN Catherine

en vue de l'obtention du diplôme de : Maîtrise ès sciences appliquées

a été dûment accepté par le jury d'examen constitué de :

M. LÉGER Pierre, Ph. D., président

M. GOULET James-A., Ph. D., membre et directeur de recherche

M. MARCOTTE Denis, Ph. D., membre

ACKNOWLEDGEMENTS

I would like to thank the Ministère des Transports, de la Mobilité durable et de l'Électrification des transports du Québec and Mathieu Lacoste for providing the data necessary for this research as well as technical assistance.

This project was funded by the Fonds de recherche du Québec en Nature et technologies (FRQNT).

REMERCIEMENTS

Un gros merci à Ha pour son support inconditionnel dans les moments plus difficiles, ainsi qu'à Ianis pour son aide précieuse. Je voudrais aussi remercier mon directeur d'études, James-A. Goulet, pour son aide, les discussions enrichissantes et sa patience.

La dernière année aurait été très différente sans vous !

Merci aussi à ma famille et amis, qui n'ont jamais arrêté de croire en moi !

RÉSUMÉ

Les infrastructures de tous les types vieillissent et leur état doit être vérifié régulièrement pour détecter les dégradations et pour planifier les opérations d'entretien. Dans les dernières décennies, les technologies des capteurs ont grandement évolué et les capteurs sont dorénavant abordables et accessibles. Ceci résulte en une utilisation grandissante des capteurs dans le domaine de la surveillance de l'état des structures, ou *Structural Health Monitoring* (SHM). SHM fait généralement référence à un système avec trois éléments principaux : le système de capteurs, le système d'analyse des données et l'évaluation de l'état de la structure.

Les modèles d'espace-état, ou *State-Space models* (SSM), sont des modèles d'estimation récursive qui évoluent avec l'arrivée de nouvelles données et sont une classe de modèle orienté données. Ils permettent d'estimer de manière dynamique les variables d'état. Les modèles bayésiens dynamiques linéaires, ou *Bayesian Dynamic Linear Models* (BDLMs), sont un type de modèle d'espace-état qui est adapté à l'inférence séquentielle.

Une difficulté associée au SHM est que des effets externes qui ne sont pas reliés à l'état réel de la structure, par exemple la température et l'humidité, sont typiquement la source d'une variation dans le comportement de la structure dont l'ordre de grandeur est comparable à l'effet d'un dommage structural significatif. Un défi dans le domaine de la surveillance de l'état des structures est de considérer les impacts de ces effets externes et de retirer ces impacts du comportement observé afin d'identifier l'état intrinsèque de la structure associé au vieillissement et à la dégradation. Quand les effets externes sont observés par des capteurs, il est possible de les inclure dans le modèle du comportement structural. Dans ce contexte, l'utilisation de multiples capteurs d'un même type a deux avantages principaux : (1) le système produit des données fiables et (2) de la variabilité spatiale est capturée par le système de capteurs. Lorsque plusieurs capteurs enregistrent le même type de données d'effet externe, inclure ces données dans un modèle n'est pas efficace en utilisant le problème de régression classique parce que les données du même type sont typiquement fortement corrélées entre elles, ce qui résulte en un problème de régression mal conditionné. Un outil commun pour s'attaquer au problème de covariables corrélées dans des analyses de régression est l'analyse par composantes principales ou *Principal Component Analysis* (PCA). Quand des variables sont liées par des corrélations linéaires, l'analyse par composantes principales permet de retirer ces corrélations en regroupant les données corrélées sur la même composante principale.

Actuellement, les BDLMs sont capables de créer des modèles qui décrivent des cas dans lesquels une observation dépendante est linéairement dépendante d'une seule observation in-

dépendante ou de plusieurs observations indépendantes mais non-corrélées. Les BDLMs ne peuvent représenter la relation qui unit un comportement structural et des multiples effets externes corrélés. Ce mémoire propose une nouvelle méthode qui est adaptée aux BDLMs et qui utilise l'analyse par composantes principales sur les observations d'effets environnementaux pour décrire la dépendance entre les effets environnementaux et la réponse structurale.

La méthode proposée consiste à représenter la dépendance linéaire entre le modèle de la réponse structurale et les impacts des modèles des effets environnementaux dans un espace transformé où les données sont non-corrélées obtenu à l'aide de l'analyse par composantes principales. Quand le nombre de composantes principales incluses dans un modèle est réduit, ceci résulte en un coût de calcul réduit et en l'élimination d'une fraction de l'information des données. Ainsi, quand la PCA est introduite dans un BDLM, la PCA est une solution pour régler le problème de système indéterminé causé par des observations corrélées utilisées comme variables indépendantes.

La méthode proposée est appliquée sur des données de SHM enregistrées sur un viaduc d'autoroute en béton armé localisé au Canada. Un jeu de données de déplacement et quatre jeux de données de température sont utilisés pour tester la nouvelle méthode. Les observations de température sont fortement corrélées sur le long terme, soit le cycle annuel, et des corrélations sont aussi visibles sur le court terme, par exemple les cycles journaliers. Quatre modèles test sont construits en utilisant la nouvelle méthode et ces modèles sont comparés avec quatre modèles d'observation du déplacement qui n'utilisent qu'une observation de température à la fois. Les conclusions de l'étude de cas sont qu'une seule composante principale est insuffisante pour représenter l'information utile contenue dans les jeux de données des multiples capteurs et que les deux dernières composantes principales contiennent de l'information qui est superflue pour prédire le déplacement de la structure. La première composante principale explique la majorité du cycle annuel de la température et les autres composantes principales regroupent l'information sur les variations de température à court terme.

Les résultats de l'étude de cas illustre que la méthode proposée réussit à traiter des effets environnementaux corrélés avec les BDLMs, ce qui était impossible avant le développement de cette méthode. De plus, l'étude de cas démontre que le fait d'inclure dans un modèle de réponse structurale les jeux de données provenant de multiples capteurs d'un même type mène à une augmentation de la capacité prédictive du modèle.

ABSTRACT

Infrastructure of all kinds are ageing over time, and their state has to be analyzed regularly to detect degradation and to plan maintenance operations. In the last decades, sensor technologies evolved significantly, and they are now affordable and accessible, which results in an increased use in the field of monitoring structural behaviours, commonly referred to as *Structural Health Monitoring* (SHM). SHM refers to a system with 3 major components : the sensor system, the data processing system and the assessment of the structural health.

A class of data-driven SHM methods employs *State-Space models* (SSM), which is a type of recursive estimation model that evolves as new data is available, and that dynamically estimates the state variables. *Bayesian Dynamic Linear Models* (BDLMs) are a class of SSM which are well suited for sequential inference.

A difficulty associated with SHM is that external conditions not related to the structure health, for example temperature and humidity, typically cause a variation in the structural behaviour that is comparable to or larger than a significant structural damage. A challenge in SHM is to take into consideration the influence of external conditions and remove their effects to extract the intrinsic response of a structure related to ageing or deterioration. When external conditions are observed, it is possible to include them in the model. The deployment of multiple sensors of the same type has two main advantages : (1) it ensures having reliable data and (2) it allows capturing spatial variability. When multiple sensors measure the same quantity, including the available data of the same type in a model cannot be resolved effectively using classical regression problem because data of the same type are typically highly correlated between each other, which results in an ill-conditioned regression problem. A common tool to tackle the issue of correlated covariates in regression analysis is *Principal Component Analysis* (PCA). When linear correlations exist between variables, PCA enables to remove those correlations by regrouping the correlated data in the same principal component.

Currently, BDLM is able to model cases where a dependent observation is linearly dependent on either a single or several independent and uncorrelated observations. It remains unable to model the relationship between structural response and multiple correlated external conditions. This master thesis proposes a new method adapted to BDLM that uses PCA on environmental effect observations in order to model the dependency between environmental effects and structural response.

The proposed method is to model the linear dependency between the structural response

model and the individual impacts of the environmental effects models transformed in an uncorrelated space obtained through PCA. Selecting a reduced number of principal components in a model allows to reduce the calculation cost and to eliminate a portion of the information from the data. Therefore, when inserted in BDLM, PCA is a solution to tackle the issue of indeterminate systems caused by correlated observations employed as independent variables.

The proposed method is applied on a case-study based on SHM data acquired on a reinforced concrete highway bridge located in Canada. One displacement dataset and four temperature datasets are used to test the proposed method. The temperature observations are strongly correlated on the long-term scale (yearly cycles), and correlations are also observable on the short-term scale (daily cycles). Four cases are built using the proposed method, and they are compared to four displacement models that use a single temperature observation at the time. The conclusions of the case-study are that one principal component is insufficient to capture the useful information from multiple sensors, and that the last two principal components contain information that is unnecessary to predict the displacement. The first principal component explains the majority of temperature's annual cycle and the others contain information about short-term variations.

The case-study results illustrate that the method is able to handle multiple correlated environmental effects in BDLM, which was not possible before, and that including in a structural response model the information from multiple environmental effect datasets leads to an increased prediction capacity.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	iii
REMERCIEMENTS	iv
RÉSUMÉ	v
ABSTRACT	vii
TABLE OF CONTENTS	ix
LIST OF TABLES	xi
LIST OF FIGURES	xii
LIST OF ABBREVIATIONS	xiv
LIST OF SYMBOLS	xv
LIST OF APPENDICES	xvii
CHAPTER 1 INTRODUCTION	1
CHAPTER 2 BAYESIAN DYNAMIC LINEAR MODELS	4
2.1 State-space Formulation	4
2.2 BDLM Components	4
2.2.1 Local Level and Local Trend Component	5
2.2.2 Periodic Component	5
2.2.3 Autoregressive Component	6
2.3 Model Assembly	6
2.4 Multiple Observations	7
2.5 Hidden States Estimation	8
2.6 Parameter Estimation	10
2.7 Limitation	11
CHAPTER 3 MODELLING ENVIRONMENTAL EFFECT DEPENDENCY USING	
PCA	12
3.1 Formulation	13

CHAPTER 4	CASE STUDY	20
4.1	Data Description	20
4.2	Description of the cases	21
4.3	Models for cases 1 to 4 using PCA	22
4.3.1	Models for Environmental Effect Observations	23
4.3.2	Model for Structural Response Observation	27
4.4	Models for cases 5 to 8 not using PCA	32
4.4.1	Model for Structural Response Observation	32
4.5	Results	34
4.6	Discussion	34
CHAPTER 5	CONCLUSION AND RECOMMENDATIONS	38
5.1	Summary of work	38
5.2	Limitations	38
5.3	Future work	39
REFERENCES		40
APPENDICES		44

LIST OF TABLES

Table 4.1	Table indicating the observations included in each case. Displacement observation is denoted \mathbf{d} and temperature observations are $\mathbf{T1}$, $\mathbf{T2}$, $\mathbf{T3}$, and $\mathbf{T4}$. m refers to the number of principal components in the model	23
Table 4.2	Log-likelihood of the temperature models over the test set	26
Table 4.3	Displacement observation log-likelihood for the studied cases using PCA	34

LIST OF FIGURES

Figure 3.1	Correlated data transformed into an uncorrelated space using PCA	12
Figure 3.2	Flowchart representing the main steps of the proposed method. The dashed box contains the elements for the method proposed in this master thesis	19
Figure 4.1	Diagrams showing the location of displacement sensor \mathbf{d} and temperature sensors T1, T2, T3, and T4 on the bridge through (a) an elevation view and (b) a cross-section. A rectangle represents a displacement sensor and a circle represents a temperature sensors	21
Figure 4.2	Raw data for the case study. The left section presents the entire dataset and the rights presents a 2-week period of data	22
Figure 4.3	Diagram of the steps with an environmental effect observations to build \mathbf{Z}^μ matrix	27
Figure 4.4	Displacement \mathbf{d} and temperature T2 raw data showing a positive correlation between displacement and temperature observations on the short term scale	29
Figure 4.5	Transformed impact of the components from the 4 temperature observations in case 4	36
Figure 4.5	Transformed impact of the components from the 4 temperature observations in case 4	37
Figure A.1	Hidden state variables estimation for the model of temperature observation T3. The left and right parts show the distribution for the entire dataset and for a period of 14 days respectively	44
Figure A.1	Hidden state variables estimation for the model of temperature observation T3. The left and right parts show the distribution for the entire dataset and for a period of 14 days respectively	45
Figure A.2	Observed data (y_t^{obs}) and model prediction for temperature T3	46
Figure B.1	Hidden state variables estimation for the displacement observation model of case 2	47
Figure B.2	Observed data (y_t^{obs}) and model prediction for displacement observation \mathbf{d} from case 2	47
Figure C.1	Hidden state variables estimation for the displacement observation model of case 8	48

Figure C.2	Observed data (y_t^{obs}) and model prediction for displacement observation	
	d from case 8	48

LIST OF ABBREVIATIONS

BDLM	Bayesian Dynamic Linear Model
HST	Hydrostatic Season Time
HSTT or HTT	Thermal Hydrostatic Season Time
HST-Grad	Gradient Hydrostatic Season Time
ICA	Independent Component Analysis
MLE	Maximum Likelihood Estimation
MLR	Multiple Linear Regression
PC	Principal Component
PCA	Principal Component Analysis
PDF	Posterior Density Function
SHM	Structural Health Monitoring
SSM	State-Space Model
SVR	Support Vector Regression

LIST OF SYMBOLS

A	Transition matrix
AR	Autoregressive component
B	Structural response
C	Observation matrix
C*	Modified observation matrix
c	Indicates the variable is grouped by observation
$\text{cov}(\cdot)$	Covariance operator
D	Dependence matrix
$d_{i,j}$	Indicates if observation i is dependent on observation j
E_i	Environmental effect i
G	Innovation covariance matrix
$i j$	i conditional on j
J	Backward Kalman gain matrix
k	Total number of hidden state variables in a model
k_j	Number of hidden state variables associated to observation j
K	Kalman gain matrix
LL	Local level component
LT	Local trend component
m	Number of principal components included in the structural response model
n	Number of environmental effect observations
P	Periodic component
P_{S1}	First sub-component of periodic component
P_{S2}	Second sub-component of periodic component
p	Refers to a given principal component
P	PCA coefficient matrix
$p(\cdot)$	Probability of \cdot
Q	Model error covariance matrix
R	Observational error covariance matrix
r	Innovation vector
s	Vector of variables transformed in the principal component space
t	Timestamps
T	Total number of timestamps in a dataset

\mathbf{v}	Vector of observational error
\mathbf{w}	Vector of model error
\mathbf{x}	Vector of hidden state variables
y	Observation
y	Number of observation in the model
\mathbf{y}	Observation vector
$\hat{\mathbf{y}}$	Prediction of observations
\mathbf{Z}	PCA's data matrix
$\mathbb{E}[\cdot]$	Expected value operator
\mathcal{P}	Vector of unknown parameters
\mathcal{P}^*	Vector of optimal parameters
Δt	Size of time step
δ	Vector of individual impacts
δ^*	Vector of individual impacts transformed in the principal component space
ε	Explained vector from PCA
μ	Expected value
σ	Standard deviation
Σ	Covariance matrix
ϕ^{AR}	Autocorrelation coefficient
$\phi^{i j}$	Regression coefficient of observation j on i
ϕ^{PC}	Scaling factor
$\boldsymbol{\phi}^{\text{PC}}$	Scaling factor matrix
ω	Angular frequency
$[\]_{\text{O}}$	Indicates in the original space
t	At time t
$t t-1$	Prior value from the Kalman filter
$t t$	Posterior value from the Kalman filter
$t T$	Posterior value from the Kalman smoother

LIST OF APPENDICES

Appendix A	ENVIRONMENTAL EFFECT OBSERVATIONS MODEL FOR T3 .	44
Appendix B	DISPLACEMENT MODEL FOR CASE 2	47
Appendix C	DISPLACEMENT MODEL FOR CASE 8	48

CHAPTER 1 INTRODUCTION

Infrastructures of all kinds are ageing over time, and their state has to be analyzed regularly to detect degradation and to plan maintenance operations. In the last decades, sensor technologies evolved significantly (Lynch and Loh, 2006), and they are now affordable and accessible, which results in an increased use in the field of monitoring structural behaviours. Similarly, methods to interpret data and to assess a structure’s health were developed and tested. This field of study, commonly referred to as *Structural Health Monitoring* (SHM), solely relies on recorded data to build empirical models which are used to assess the health of a structure. SHM refers to a system with of 3 major elements : the sensor system, the data processing system and the assessment of the structural health, as defined by Ni et al. (2005). In this thesis, we focus on data-driven SHM. One major advantage of data-driven techniques is that they are generic because they do not require previous knowledge about the mechanical and/or geometric properties of the structure to predict the behaviour.

One such data-driven technique employs regression methods. Regression methods enable to build a model defined by the equation $\mathbf{y} = \mathbf{g}(\mathbf{x})$, where \mathbf{y} is a structural response observation that is composed of functions $\mathbf{g}(\cdot)$ of time varying covariates \mathbf{x} . The time-dependent covariates can be observed or unobserved (hidden) variables. In this kind of method, the model is built and optimized using the training set, but the model remains unchanged when new data is added. An example of Regression method is *Multiple Linear Regression* (MLR) models, which are based on the principle that some measured quantities, such as temperature, help describe a structure’s response through linear regressions. For instance, Tatin et al. (2015) and Léger and Leclerc (2007) use polynomial functions of the water height, among other functions and measured quantities, to describe the structural response of a dam. Artificial Neural Network (Mata, 2011; Xia et al., 2011a) has also shown good performances in SHM.

Another class of SHM method employs the *State-Space model* (SSM), which is a type of recursive estimation model that evolves as new data is available. This type of model has the advantage of dynamically estimating hidden states, which means that current hidden states depend on the value from previous time step, as opposed to Regression methods. Examples of common SSMs methods are the Autoregressive models (Peeters and De Roeck, 2001) and the *Bayesian Dynamic Linear Models* (BDLM) (Goulet, 2017), which use linear regression to model dependencies between hidden covariates. BDLMs are a class of state-space models (Särkkä, 2013) which are well suited for sequential inference (Goulet, 2017). BDLMs are able to decompose observed time series into a set of hidden state variables. One major advantage

of BDLMs is that the behaviour of the extracted hidden state variables can vary over time. Recent applications have demonstrated the potential of BDLMs to track time-varying baseline response from real datasets and detect anomalies (Nguyen and Goulet, 2018).

A difficulty associated with SHM is that external conditions not related to the health of the structure usually strongly affect the behaviour. External conditions are commonly separated in two categories: operational conditions and environmental effects (Ni et al., 2005). Operational conditions are loads that originate from the normal operation of a structure (traffic load, human occupation). Environmental effects are defined as loads that are generated from a variation in the environmental conditions of a structure (temperature, wind, water level, etc.). Many authors agree that a variation in an external condition does have a considerable impact on the structural response (Li et al., 2016; Hua et al., 2007; Westgate et al., 2014; Limongelli et al., 2016; Peeters et al., 2000; Mata et al., 2013; Yuen and Kuok, 2010a). External condition typically cause a variation in the structural behaviour that is comparable to or larger than a significant structural damage, as explained by Xia et al. (2012), Yuen and Kuok (2010b), and Ni et al. (2005).

A main challenge in SHM is to take into consideration the influence of external conditions and remove their effects to extract the intrinsic response of a structure (i.e. baseline response). In that context, the observations related to the intrinsic structural response i.e., displacement or frequency, are considered as dependent variables that are affected by independent variables, which are covariates related to the external conditions. When the external conditions are not observed (i.e., hidden covariates), one possibility is to model them using periodic functions (Nguyen and Goulet, 2018; Tatin et al., 2015; Léger and Leclerc, 2007). On the opposite, when external conditions are observed, it is possible to include them in the model as observed covariates. Because sensors are now widely available, structures having multiple sensors measuring the same quantity are frequent (Xia et al., 2012; Ni et al., 2005; Limongelli et al., 2016). For instance, this situation occurs when several sensors measuring the same physical quantity are positioned at different locations. Moreover, the deployment of multiple sensors of the same type are common in SHM for two other reasons : (1) it ensures having reliable data, as a sensor could fail or drift over time, (2) it allows capturing spatial variability. In the case of environmental effect sensors, the spatial variability can have a significant impact on the material properties and consequently on the structural response, as Enckell et al. (2011) mentionned. Xia et al. (2011b) showed that catching the spatial distribution is important to obtain a better quality of predictions. However, when this kind of data is available, one may wonder whether there is an advantage to include all the available data to improve the model or not. Including in a model all (or a portion of) the available data of the same type cannot be resolved effectively using classical regression problem because data of the same type are

typically highly correlated between each other, which results in an ill-conditioned regression problem.

A common tool to tackle the issue of correlated covariates in regression analysis is *Principal Component Analysis* (PCA). PCA is a space transformation originally developed by Pearson (1901). It converts a set of correlated variables into a new set of linearly uncorrelated variables using linear transformation (Abdi and Williams, 2010). One of its early usage in the context of SHM was made by Worden et al. (2000). Many other authors followed and also exploited the advantages of PCA. For instance, Yan et al. (2005) and Magalhães et al. (2012) developed different methods where PCA is used to remove the effect of unobserved environmental effects from a structural response observation. Malekzadeh et al. (2015) performed PCA on a moving window of data to extract damage indexes from the first principal component. Hua et al. (2007) used PCA in a different way, and developed a method that extracts useful information from correlated environmental effect observation using PCA to model the structural response with the *Support Vector Regression* (SVR) technique. They demonstrated that PCA-compressed data are more effective at describing a structural response than correlated data, and they improve the prediction capacity of the model.

The current BDLM approach enables to model the dependency with uncorrelated variables using linear regression. However, it does not handle correlated environmental effects time series. This master thesis proposes a new method adapted to BDLM that uses PCA as a dimension reducing technique for environmental effect observations in order to model the dependency between environmental effects and the structural response. First, Chapter 2 details the theory for Bayesian Dynamic Linear Models. Second, the proposed methodology to model environmental effects dependencies using PCA is detailed in Chapter 3. Then, Chapter 4 presents an application with real data from a bridge located in Canada. Finally, Chapter 5 concludes with a summary of work, the limitations and future improvements.

CHAPTER 2 BAYESIAN DYNAMIC LINEAR MODELS

This chapter presents the details about BDLMs as presented by West and Harrison (1999) and Goulet (2017). The general equations for building a model and learning its parameters are presented and explained.

2.1 State-space Formulation

This section presents BDLMs, which is a special class of SSMs. A BDLM is defined by an observation model and a transition model. The observation model, employed to describe the relation between observations $\mathbf{y}_t = [y_1, y_2, \dots, y_y]^\top$ and hidden state variables $\mathbf{x}_t = [x_1, x_2, \dots, x_k]^\top$, where k is the total number of hidden state variables, is described by

$$\mathbf{y}_t = \mathbf{C}_t \mathbf{x}_t + \mathbf{v}_t, \quad \begin{cases} \mathbf{y}_t \sim \mathcal{N}(\mathbb{E}[\mathbf{y}_t], \text{cov}[\mathbf{y}_t]) \\ \mathbf{x}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) \\ \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t), \end{cases} \quad (2.1)$$

where \mathbf{C}_t is the observation matrix, the hidden state variables \mathbf{x}_t all follow a Gaussian distribution with mean $\boldsymbol{\mu}_t$ and covariance matrix $\boldsymbol{\Sigma}_t$, and \mathbf{v}_t is the Gaussian measurement error with zero mean and covariance matrix \mathbf{R}_t . The *hidden state variables* \mathbf{x}_t are components of the observations \mathbf{y}_t that can not be directly observed. The transition model describing the dynamics of the hidden state variables \mathbf{x}_t over time, is defined as

$$\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_t), \quad (2.2)$$

where \mathbf{A}_t is the transition matrix, and \mathbf{w}_t represents the Gaussian model errors with zero mean and covariance matrix \mathbf{Q}_t . To simplify the notation, the subscript t is dropped for the model matrices $\{\mathbf{A}_t, \mathbf{C}_t, \mathbf{Q}_t, \mathbf{R}_t\}$ even if in practice model matrices may vary over time.

2.2 BDLM Components

The general idea behind BDLM consists in decomposing an observation into a vector of hidden state variables. Using Equations 2.1 and 2.2, a prediction of the distribution of observation(s) \mathbf{y}_t is estimated with data from the previous timesteps, and the estimation can be updated with incoming data. By definition, a hidden state variable is one that is not directly observed. The hidden state variables can be a *local level* component, a *trend*,

a *periodic* component, and an *autoregressive* component, where each hidden state variable serves different roles in the BDLMs.

2.2.1 Local Level and Local Trend Component

The local level component is employed to model, for example, the structural behaviour baseline without external conditions. It also represents the general baseline in the case of an external condition observation. The corresponding hidden state variable and model matrices are

$$\begin{aligned} \mathbf{x}^{\text{LL}} &= x^{\text{LL}}, \\ \mathbf{A}^{\text{LL}} &= \begin{bmatrix} 1 \end{bmatrix}, \\ \mathbf{C}^{\text{LL}} &= \begin{bmatrix} 1 \end{bmatrix}, \\ \mathbf{Q}^{\text{LL}} &= (\sigma^{\text{LL}})^2. \end{aligned} \tag{2.3}$$

When necessary, the local trend is used to describe the rate of change in the baseline component that is, the variation in the baseline over time. The corresponding hidden state variable and model matrices are

$$\begin{aligned} \mathbf{x}^{\text{LT}} &= \begin{bmatrix} x^{\text{LL}} \\ x^{\text{LT}} \end{bmatrix}, \\ \mathbf{A}^{\text{LT}} &= \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}, \\ \mathbf{C}^{\text{LT}} &= \begin{bmatrix} 1 & 0 \end{bmatrix}, \\ \mathbf{Q}^{\text{LT}} &= (\sigma^{\text{LT}})^2 \begin{bmatrix} \frac{\Delta t^3}{3} & \frac{\Delta t^2}{2} \\ \frac{\Delta t^2}{2} & \Delta t \end{bmatrix}, \end{aligned} \tag{2.4}$$

where Δt is the size of the time step. When the local trend varies over time, it is possible to model this variation with a *local acceleration*, and with higher order models such as a *local jerk*.

2.2.2 Periodic Component

The periodic component describes harmonic cyclic variations such as daily cycles of temperature. For each periodic component, there are 2 associated hidden state variables. Only the first hidden state variable contributes to the observation y_t (Goulet, 2017; West and Harrison, 1999). The hidden state variables and model matrices follow

$$\begin{aligned}
\mathbf{x}^P &= \begin{bmatrix} x^{P,S1} \\ x^{P,S2} \end{bmatrix}, \\
\mathbf{A}^P &= \begin{bmatrix} \cos(\omega^P \Delta t) & \sin(\omega^P \Delta t) \\ -\sin(\omega^P \Delta t) & \cos(\omega^P \Delta t) \end{bmatrix}, \\
\mathbf{C}^P &= \begin{bmatrix} 1 & 0 \end{bmatrix}, \\
\mathbf{Q}^P &= \begin{bmatrix} (\sigma^{P,S1})^2 & 0 \\ 0 & (\sigma^{P,S2})^2 \end{bmatrix},
\end{aligned} \tag{2.5}$$

where $\omega^P = \frac{2\pi}{P}$ is the angular frequency.

A model may contain several periodic components with different periods.

2.2.3 Autoregressive Component

The autoregressive component is used to capture the time-dependent residual between the model prediction and the observation at each time step. If a pattern is recognizable in the autoregressive term, it means that a component might be missing in the model. The AR component used in BDLM is built as follows

$$\begin{aligned}
\mathbf{x}^{\text{AR}} &= x^{\text{AR}}, \\
\mathbf{A}^{\text{AR}} &= \begin{bmatrix} \phi^{\text{AR}} \end{bmatrix}, \\
\mathbf{C}^{\text{AR}} &= \begin{bmatrix} 1 \end{bmatrix}, \\
\mathbf{Q}^{\text{AR}} &= (\sigma^{\text{AR}})^2.
\end{aligned} \tag{2.6}$$

In order to have a stationary process, the value of the autocorrelation coefficient ϕ^{AR} has to be between 0 and 1.

2.3 Model Assembly

To form the model, the components are assembled and they build the model matrices. For the hidden state variable vector \mathbf{x}_t , the elements from the components are assembled to form a column vector such that, for instance, an observation model with a local level and a periodic component with period P1 forms the vector

$$\mathbf{x}_t = \begin{bmatrix} \mathbf{x}^{\text{LL}} \\ \mathbf{x}^{\text{P1}} \end{bmatrix} = \begin{bmatrix} x^{\text{LL}} \\ x^{\text{P1},S1} \\ x^{\text{P1},S1} \end{bmatrix}. \tag{2.7}$$

For the transition matrix \mathbf{A} and the model error covariance matrix \mathbf{Q} , the matrices from the components are assembled in a diagonal matrix, that is, for the previously mentioned example

$$\begin{aligned}\mathbf{A} &= \text{block diag.}(\mathbf{A}^{\text{LL}}, \mathbf{A}^{\text{P1}}) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(\omega^{\text{P1}}\Delta t) & \sin(\omega^{\text{P1}}\Delta t) \\ 0 & -\sin(\omega^{\text{P1}}\Delta t) & \cos(\omega^{\text{P1}}\Delta t) \end{bmatrix} \\ \mathbf{Q} &= \text{block diag.}(\mathbf{Q}^{\text{LL}}, \mathbf{Q}^{\text{P1}}) = \begin{bmatrix} (\sigma^{\text{LL}})^2 & 0 & 0 \\ 0 & (\sigma^{\text{P},\text{S1}})^2 & 0 \\ 0 & 0 & (\sigma^{\text{P},\text{S2}})^2 \end{bmatrix}.\end{aligned}\tag{2.8}$$

In the observation matrix, the matrices from the components are concatenated to form a line vector

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}^{\text{LL}} & \mathbf{C}^{\text{P1}} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \end{bmatrix}.\tag{2.9}$$

2.4 Multiple Observations

BDLM can be employed to model several observations simultaneously. When multiple observations are employed within a same model ($\mathbf{y}_t = [y_1, y_2, \dots, y_j, \dots, y_y]^\top$), the hidden state variables associated with each observation are grouped in a hidden state vector $\mathbf{x}_t = [\mathbf{x}_1^c, \mathbf{x}_2^c, \dots, \mathbf{x}_j^c, \dots, \mathbf{x}_y^c]^\top$ where $\mathbf{x}_j^c \in \mathbb{R}^{k_j}$, k_j being the number of hidden state variables associated with the observation y_j . The model matrices are built using the same method as presented in Section 2.3. In practical applications, it is common to have dependencies between observations, e.g. the structural response (first observation) depends on the air temperature (second observation). Linear dependencies between observations are defined in the observation matrix \mathbf{C} . The component-wise representation of the observation matrix is

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}^{c,1} & \mathbf{C}^{c,1|2} & \dots & \mathbf{C}^{c,1|j} & \dots & \mathbf{C}^{c,1|y} \\ \mathbf{C}^{c,2|1} & \mathbf{C}^{c,2} & \dots & \mathbf{C}^{c,2|j} & \dots & \mathbf{C}^{c,2|y} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{C}^{c,i|1} & \mathbf{C}^{c,i|2} & \dots & \mathbf{C}^{c,i|j} & \dots & \mathbf{C}^{c,i|y} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{C}^{c,y|1} & \mathbf{C}^{c,y|2} & \dots & \mathbf{C}^{c,y|j} & \dots & \mathbf{C}^{c,y} \end{bmatrix},\tag{2.10}$$

where $\mathbf{C}^{c,i|j}$ are row vectors with k_j elements. In the special case where there is no dependence between observations, only components on the diagonal of \mathbf{C} are non-zero. The existence of dependence between observations is encoded by binary variables $d_{i,j} \in \{0, 1\}$ that are grouped

in a dependence matrix \mathbf{D} where when $d_{i,j} = 0$, $\mathbf{C}^{c,i|j} = [\mathbf{0}]$ and when $d_{i,j} = 1$, $\mathbf{C}^{c,i|j} = [\phi_1^{i|j}, \phi_2^{i|j}, \dots, \phi_{k_j}^{i|j}] \in \mathbb{R}^{k_j} \forall i \neq j$. The variable $\phi_k^{i|j}$ is a regression coefficient describing the dependency between observation i and the k^{th} hidden variable of observation j . Regression coefficients $\phi_k^{i|j}$ needs to be estimated using the *Maximum Likelihood Estimation* presented in Section 2.6. Note that $d_{i,j} = 1 \forall i = j$ because an observation is correlated to itself. At a given time t , when $d_{i,j} = 1 \forall i \neq j$, the variation in dependent observation i caused by an independent observation j , or the impact $\delta_t^{i|j}$ of observation j on observation i , is quantified by

$$\delta_t^{i|j} = \mathbf{C}^{c,i|j}[\mathbf{x}_j^c]_t. \quad (2.11)$$

The BDLM is able to distinguish the long term from short term impacts of variation in the environmental effects on a structural response, which is done with the different regression coefficients in the observation matrix $\mathbf{C}^{c,i|j} = [\phi_1^{i|j}, \phi_2^{i|j}, \dots, \phi_{k_j}^{i|j}]$. For example, long term variation in an environmental effect, such as temperature, might generate a negligible thermal gradient in a structure because it occurs over a long period of time, as opposed to short term variation, due to thermal inertia. For this reason, short term environmental effect variation induces different internal efforts than long term variation that could lead to a change in the structural response. Modelling this phenomena increases the flexibility of the model and offers a better prediction of the displacement data.

2.5 Hidden States Estimation

The posterior *probability density function* (pdf) of hidden state variables \mathbf{x}_t at $t \in [0, T]$ given the observations $\mathbf{y}_{1:t}$ are calculated using the Kalman filter algorithm (Murphy, 2012; Welch and Bishop, 2001). The Kalman filter is an iterative process involving two steps: the prediction step and the measurement step. The *prediction step* enables to obtain the prior distribution of the hidden state variables \mathbf{x}_t using the knowledge of x_{t-1}

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) &= \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) && \text{Prior state estimation} \\ \boldsymbol{\mu}_{t|t-1} &\triangleq \mathbf{A}_t \boldsymbol{\mu}_{t-1|t-1} && \text{Prior expected value} \\ \boldsymbol{\Sigma}_{t|t-1} &\triangleq \mathbf{A}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{A}_t^\top + \mathbf{Q}_t && \text{Prior covariance} \end{aligned} \quad (2.12)$$

The *measurement step* enables to estimate the posterior distribution of hidden state variables \mathbf{x}_t using the observation(s) at time t .

$$\begin{aligned}
p(\mathbf{x}_t | \mathbf{y}_{1:t}) &= \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) && \text{Posterior state estimation} \\
\boldsymbol{\mu}_{t|t} &= \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t \mathbf{r}_t && \text{Posterior expected value} \\
\boldsymbol{\Sigma}_{t|t} &= (\mathbf{I} - \mathbf{K}_t \mathbf{C}_t) \boldsymbol{\Sigma}_{t|t-1} && \text{Posterior covariance} \\
\mathbf{r}_t &\triangleq \mathbf{y}_t - \hat{\mathbf{y}}_t && \text{Innovation vector} \\
\hat{\mathbf{y}}_t &\triangleq \mathbb{E}[\mathbf{y}_t | \mathbf{y}_{1:t-1}] = \mathbf{C}_t \boldsymbol{\mu}_{t|t-1} && \text{Predicted observations vector} \\
\mathbf{K}_t &\triangleq \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}_t^\top \mathbf{G}_t^{-1} && \text{Kalman gain matrix} \\
\mathbf{G}_t &\triangleq \mathbf{C}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}_t^\top \mathbf{G}_t^{-1} && \text{Innovation covariance matrix.}
\end{aligned} \tag{2.13}$$

where

$$\boldsymbol{\mu}_{t|t} = \mathbb{E}[\mathbf{x}_t | \mathbf{y}_{1:t}]$$

$$\boldsymbol{\Sigma}_{t|t} = \text{cov}[\mathbf{x}_t | \mathbf{y}_{1:t}].$$

$\boldsymbol{\mu}_{t|t}$ and $\boldsymbol{\Sigma}_{t|t}$ refer respectively to the posterior expected value and the posterior covariance matrix of \mathbf{x}_t at time t , given the observations $\mathbf{y}_{1:t}$, that is, $(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$ are the output of the Kalman filter. For the posterior expected value, the Kalman gain serves as a way to weight the information from the new observations \mathbf{y}_t compared to the information from prior knowledge. The short form of the filtering step for estimating hidden state at a time t from hidden state estimates at $t - 1$ and the observation \mathbf{y} at t is

$$(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) = \text{Filter}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{y}_t, \mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}). \tag{2.14}$$

The offline estimation for the hidden state variables \mathbf{x}_t at time t , that is, the estimation of hidden state variables at time t using the entire set of available data $(\mathbf{y}_{1:T})$, can be performed using the Kalman smoother (Murphy, 2012). Similarly to the Kalman filter, the Kalman smoother is employed to estimate the posterior distribution for the hidden state variables \mathbf{x}_t that is also assumed to be a multivariate Gaussian distribution following

$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}). \tag{2.15}$$

The difference to the hidden state variable distributions estimated by Kalman filter is that the smoothed hidden state variables are conditioned to the observations $\mathbf{y}_{1:T}$, while the filtered hidden states are conditional only to the observations from the previous time steps $\mathbf{y}_{1:t}$. In addition, the initial values for the hidden state variables estimated using Kalman smoother are those obtained from the last step of the Kalman filter i.e., $(\boldsymbol{\mu}_{T|T}, \boldsymbol{\Sigma}_{T|T})$, making the

Kalman smoother a backward process. The equations for the Kalman smoother follows

$$\begin{aligned}
p(\mathbf{x}_t | \mathbf{y}_{1:T}) &= \mathcal{N}(\mathbf{x}_T | \boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) \\
\boldsymbol{\mu}_{t|T} &= \boldsymbol{\mu}_{t|t} + \mathbf{J}_t(\boldsymbol{\mu}_{t+1|T} - \boldsymbol{\mu}_{t+1|t}) && \text{Posterior expected value} \\
\boldsymbol{\Sigma}_{t|T} &= \boldsymbol{\Sigma}_{t|t} + \mathbf{J}_t(\boldsymbol{\Sigma}_{t+1|T} - \boldsymbol{\Sigma}_{t+1|t})\mathbf{J}_t^\top && \text{Posterior covariance} \\
\mathbf{J}_t &\triangleq \boldsymbol{\Sigma}_{t|t}\mathbf{A}_{t+1}^\top\boldsymbol{\Sigma}_{t+1|t}^{-1} && \text{Backward Kalman gain matrix}
\end{aligned} \tag{2.16}$$

The Kalman smoother enables to estimate initial values for the hidden state variables. The smoothing step for estimating hidden state variables \mathbf{x}_t using $\mathbf{y}_{1:T}$ is summarized in the short form following

$$(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) = \text{Smoother}(\boldsymbol{\mu}_{t+1|T}, \boldsymbol{\Sigma}_{t+1|T}, \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \mathbf{A}, \mathbf{Q}). \tag{2.17}$$

2.6 Parameter Estimation

The model matrices $\{\mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}\}$ contain several unknown model parameters \mathcal{P} to be inferred from data. For this purpose, the chosen approach is the *Maximum Likelihood Estimation* (MLE). The MLE consists in identifying the optimal model parameters \mathcal{P}^* by maximizing the likelihood function

$$p(\mathbf{y}_{1:T} | \mathcal{P}) = \prod_{t=1}^T p(\mathbf{y}_t | \mathcal{P}), \tag{2.18}$$

representing the likelihood of observations \mathbf{y}_t conditional to the values of the parameters \mathcal{P} , for $t = \{1, 2, \dots, T\}$. Equation 2.18 is based on the hypothesis that observations $\mathbf{y}_{1:T}$ are independent, and it is valid if and only if this hypothesis is respected. Because the total number of timestamps T is generally large and for numerical stability, the logarithm of the likelihood is calculated

$$\ln p(\mathbf{y}_{1:T} | \mathcal{P}) = \sum_{t=1}^T \ln p(\mathbf{y}_t | \mathcal{P}). \tag{2.19}$$

In BDLMs, the marginal likelihood $p(\mathbf{y}_t | \mathcal{P})$ is a multivariate Gaussian distribution so that Equation 2.19 can be expanded as

$$\ln p(\mathbf{y}_{1:T} | \mathcal{P}) = \sum_{t=1}^T \ln \left[\mathcal{N}(\mathbf{y}_t; \mathbf{C}\boldsymbol{\mu}_{t|t-1}, \mathbf{R} + \mathbf{C}\boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^\top) \right], \tag{2.20}$$

where the prior mean value $\boldsymbol{\mu}_{t|t-1}$ and prior covariance matrix $\boldsymbol{\Sigma}_{t|t-1}$ for the hidden state variables are obtained using the transition model in Equation 2.12. The optimization task is carried out using optimization methods such as the *Newton-Raphson* algorithm (Gelman

et al., 2014). The Newton-Raphson algorithm is an iterative process that linearizes a function at a given point to obtain a better estimate of a root at each iteration. In this case, the Newton-Raphson algorithm is used to estimate the optimized model parameters \mathcal{P}^* by maximizing the joint prior probability of observations $\mathbf{y}_{1:T}$.

2.7 Limitation

Currently, BDLM is able to model cases where a dependent observation is linearly dependent on either a single or several independent and uncorrelated observations with the regression coefficients $\phi_k^{i|j}$. BDLM was successfully applied to model the natural frequency of a bridge which depends on the temperature and traffic load by Goulet and Koo (2018). Nevertheless, it remains unable to model the relationship between structural responses and multiple correlated external conditions. This limitation occurs because performing regression on linearly dependent variables is intrinsically underdetermined i.e., there is an infinite number of possible solutions. The next chapter presents the method proposed in this thesis to tackle this limitation by including a principal component decomposition in the existing BDLM formulation.

CHAPTER 3 MODELLING ENVIRONMENTAL EFFECT DEPENDENCY USING PCA

Principal Component Analysis (PCA) enables to transform the data from the original space to a set of orthogonal axes named the principal components (PCs). Each PC explains a decreasing quantity of the data's variance (Jolliffe, 2002; Abdi and Williams, 2010). Note that a PC does not necessarily have a physical meaning because of the transformation. When linear correlations exist between variables, PCA enables to remove those correlations by grouping the correlated data in a same PC. For instance, the data points presented in Figure 3.1(a) are correlated i.e., when x increases, so does y . Using PCA, the data is transformed in the PC space, shown in Figure 3.1(b). In this last figure, the data points are uncorrelated, that is, the value in axis PC 1 does not give information about the value in axis PC 2. Note that the span of values for PC 1 is larger than for PC 2 because PC 1 explains a majority of the variance in the data. Therefore, the data could be compressed to PC 1 with a limited information loss.

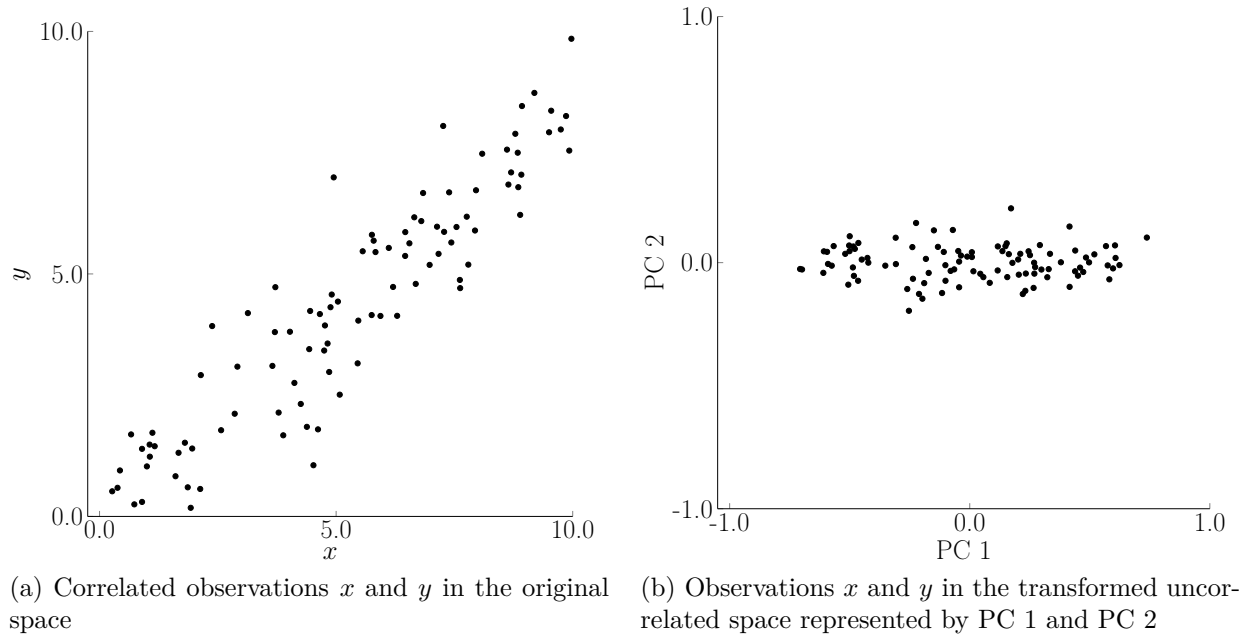


Figure 3.1 Correlated data transformed into an uncorrelated space using PCA

The PCA transformation is commonly used as a dimension reduction technique as it allows keeping a maximum amount of information within a smaller number of transformed variables.

The information contained in the combination of all PCs is equal to the information in the original space. In other words, PCA enables to transform a dataset from a correlated space to an uncorrelated space without information loss, which is an essential characteristic for this research. As explained in Chapter 1, datasets in SHM are commonly composed of correlated time series. In this context, PCA provides the means to perform linear regression between a dependent variable (e.g. structural response) and several sensors of a given type of environmental effect (e.g. temperature).

The proposed method is to model the linear dependency between the structural response model and the environmental effects models transformed in an uncorrelated space. Making this task compatible with the BDLM method presented in Chapter 2 requires a new formulation for the observation matrix \mathbf{C} . In BDLM, environmental effects are decomposed in a vector of hidden state variables such as a local level, daily and seasonal periodic components, and an auto-regressive component. These components are then transformed using PCA from the original space to the orthogonal uncorrelated space. The effect of a PC on the structural response is scaled by a regression coefficient ϕ_p^{PC} where p refers to the principal component number. Each environmental effect observation contributes to each one of the principal components which themselves are linear combinations of the data. The transformation of the independent observations in the PC space is detailed in the following section. When inserted in BDLM, PCA is a solution to tackle the issue of indeterminate systems caused by correlated observations employed as independent variables.

3.1 Formulation

We define a dependent observation y_t^{B} describing, for example, a structural response, and a set of n independent observations \mathbf{y}_t^{E} describing, for example, environmental effects. The observation vector for the model is

$$\begin{aligned} \mathbf{y}_t &= \begin{bmatrix} y_t^{\text{B}} \\ \mathbf{y}_t^{\text{E}} \end{bmatrix}, \\ \mathbf{y}_t^{\text{E}} &= [y_t^{\text{E}_1}, y_t^{\text{E}_2}, \dots, y_t^{\text{E}_i}, \dots, y_t^{\text{E}_n}]^{\text{T}}. \end{aligned} \quad (3.1)$$

A model per environmental effect observation $y_t^{\text{E}_i}$ is built in order to decompose each observation into subcomponents. The observation and transition models for each E_i observation are respectively

$$y_t^{\text{E}_i} = \mathbf{C}^{\text{E}_i} \mathbf{x}_t^{\text{E}_i} + v_t^{\text{E}_i}, \quad v_t^{\text{E}_i} \sim \mathcal{N}(0, \sigma_v^{\text{E}_i}), \quad (3.2)$$

$$\mathbf{x}_t^{\mathbf{E}_i} = \mathbf{A}^{\mathbf{E}_i} \mathbf{x}_{t-1}^{\mathbf{E}_i} + \mathbf{w}_t^{\mathbf{E}_i}, \quad \mathbf{w}_t^{\mathbf{E}_i} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{\mathbf{E}_i}). \quad (3.3)$$

The model parameters involved in the definition of the model matrices $\{\mathbf{A}^{\mathbf{E}_i}, \mathbf{C}^{\mathbf{E}_i}, \mathbf{Q}^{\mathbf{E}_i}, \sigma_v^{\mathbf{E}_i}\}$ are estimated using the MLE method presented in Section 2.6. Using the Kalman filter and smoother from Equations 2.14 and 2.17, the hidden state variables $\mathbf{x}^{\mathbf{E}_i}$ for each environmental effect observation are estimated by a mean vector and covariance matrix at each timestamp

$$\mathbf{x}_t^{\mathbf{E}_i} \sim \mathcal{N}(\boldsymbol{\mu}_{t|T}^{\mathbf{E}_i}, \boldsymbol{\Sigma}_{t|T}^{\mathbf{E}_i}). \quad (3.4)$$

Before processing PCA, the data matrix \mathbf{Z} is built from the hidden state variables of the environmental effects that cause a variation in the structural response. The equation defining environmental effects variation matrix \mathbf{Z} at time t for the i^{th} environmental effect observation is

$$[\mathbf{Z}]_{t,i} = z_{t,i} = \mathbf{C}^{\mathbf{E}_i*} \mathbf{x}_t^{\mathbf{E}_i}, \quad (3.5)$$

where $\mathbf{C}^{\mathbf{E}_i*}$ is equal to $\mathbf{C}^{\mathbf{E}_i}$ with the difference that $\mathbf{C}^{\mathbf{E}_i, \text{LL}} = [0]$ for the local level component. Each term in \mathbf{Z} corresponds to a sum of the hidden state variables from a given environmental effect that have an impact on the dependent variable $y_t^{\mathbf{B}}$. The local level component is excluded because it is constant over time, and consequently, does not have an impact on the dependent variable. Data variation matrix $[\mathbf{Z}]_{T \times n}$ has n columns corresponding to the n environmental effect observations and each line corresponds to a different timestamp $t \in [1, 2, \dots, T]$. Removing the local level component ensures that each column of \mathbf{Z} has a zero-mean (Shlens, 2014). PCA decomposition can be performed on scalars and not distributions like those describing the hidden state variables. To overcome this problem, the data matrix employed to compute the PCA decomposition is built using only the posterior expected value of hidden state variables

$$[\mathbf{Z}^\mu]_{t,i} = z_{t,i}^\mu = \mathbf{C}^{\mathbf{E}_i*} \boldsymbol{\mu}_{t|T}^{\mathbf{E}_i}. \quad (3.6)$$

PCA's data matrix \mathbf{Z}^μ dimensions are $[T \times n]$, and its composition is

$$\mathbf{Z}^\mu = \begin{bmatrix} z_{1,1}^\mu & z_{1,2}^\mu & \cdots & z_{1,i}^\mu & \cdots & z_{1,n}^\mu \\ z_{2,1}^\mu & z_{2,2}^\mu & \cdots & z_{2,i}^\mu & \cdots & z_{2,n}^\mu \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ z_{t,1}^\mu & z_{t,2}^\mu & \cdots & z_{t,i}^\mu & \cdots & z_{t,n}^\mu \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ z_{T,1}^\mu & z_{T,2}^\mu & \cdots & z_{T,i}^\mu & \cdots & z_{T,n}^\mu \end{bmatrix},$$

where the matrix \mathbf{Z}^μ is the data matrix used to compute PCA decomposition. PCA decomposition on matrix \mathbf{Z}^μ enables to obtain the PCA coefficient matrix \mathbf{P} .

$$\mathbf{P} = \begin{bmatrix} P_{1,1} & P_{1,2} & \cdots & P_{1,p} & \cdots & P_{1,n} \\ P_{2,1} & P_{2,2} & \cdots & P_{2,p} & \cdots & P_{2,n} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ P_{i,1} & P_{i,2} & \cdots & P_{i,p} & \cdots & P_{i,n} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ P_{n,1} & P_{n,2} & \cdots & P_{n,p} & \cdots & P_{n,n} \end{bmatrix},$$

where $i \in \{1, 2, \dots, n\}$ refers to a given environmental effect and $p \in \{1, 2, \dots, n\}$ refers to a given PC. The dimensions of matrix \mathbf{P} are $[n \times n]$ where the p^{th} column corresponds to the coefficients to transform the original data at time t ($[\mathbf{Z}^\mu]_{t,1:n}$) into the p^{th} PC. PCA decomposition also returns the explained vector $\boldsymbol{\varepsilon}$, whose size is $[n \times 1]$. The p^{th} value in this vector corresponds to the percentage of the total variance of the original data (\mathbf{Z}^μ) that is explained by the p^{th} PC. Those values are employed to (1) quantify the variance explained by each PC, and (2) identify the possibility to model the data in a transformed space with a reduced number of dimensions. Note that using the posterior expected value of \mathbf{x}_t from the Kalman smoother, which is $\boldsymbol{\mu}_{t|T}$, is a simplification to process PCA analysis in a computationally reasonable way. In order to transform the environmental effects in the PC space, the entire set of environmental effects has to be treated at once. At time t , the PCA coefficient matrix \mathbf{P} is used to obtain the transformed variation of environmental effects in the PC space \mathbf{s}_t

$$\begin{aligned} \mathbf{s}_t &= \mathbf{P}^\top \mathbf{z}_t^\top \\ &= \mathbf{P}^\top \left(\mathbf{C}^{\text{E}*} \mathbf{x}_t^{\text{E}} \right), \end{aligned} \tag{3.7}$$

where

$$\mathbf{z}_t = [z_{t,E_1}, z_{t,E_2}, \dots, z_{t,E_i}, \dots, z_{t,E_n}]$$

$$\mathbf{C}^{\mathbf{E}^*} = \text{block diag}(\mathbf{C}^{\mathbf{E}_i^*}), i = \{1, 2, \dots, n\} = \begin{bmatrix} \mathbf{C}^{\mathbf{E}_1^*} & 0 & \dots & 0 & \dots & 0 \\ 0 & \mathbf{C}^{\mathbf{E}_2^*} & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{C}^{\mathbf{E}_i^*} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \dots & \mathbf{C}^{\mathbf{E}_n^*} \end{bmatrix}$$

$$\mathbf{x}_t^{\mathbf{E}} = \begin{bmatrix} \mathbf{x}_t^{\mathbf{E}_1} \\ \mathbf{x}_t^{\mathbf{E}_2} \\ \vdots \\ \mathbf{x}_t^{\mathbf{E}_i} \\ \vdots \\ \mathbf{x}_t^{\mathbf{E}_n} \end{bmatrix}.$$

Each value in the vector \mathbf{s}_t represents the variation of the entire set of environmental effects observations at time t in a given PC. Each principal component p has its own scaling factor ϕ_p^{PC} , $p \in \{1, 2, \dots, n\}$. These scaling factors are treated as unknown parameters to be inferred from data. They allow to weight the impact carried by a given principal component. The scaling factors are integrated in BDLM by multiplying the transformed data in the p^{th} principal component ($\mathbf{s}_{t,p}$) by ϕ_p^{PC} . Note that the percentage of variance explained by a given PC in the explained vector $\boldsymbol{\varepsilon}$ is a different concept than the impact of that given PC on a dependent observation scaled with the scaling factor ϕ_p^{PC} . For instance, the first and second PC of a dataset always explain a portion of the total variance with $\varepsilon_1 > \varepsilon_2$, but the contribution of PC 2 could be larger than that of PC 1 i.e., $|\phi_1^{\text{PC}}| < |\phi_2^{\text{PC}}|$.

In the special case where the number of PCs m to be included in the dependent observation model is smaller than the total number of environmental effect observations n , the number of regression factors ϕ_p^{PC} to be calibrated is reduced because the corresponding regression factors are forced to be equal to zero. Reducing the number of PCs reduces the number of unknown parameters as it eliminates the need for the estimation of regression factor(s) ϕ_p^{PC} , $\forall p > m$. This allows for reduced calculation cost with limited information loss, which depends on the values in $\boldsymbol{\varepsilon}$.

To ensure that some flexibility remains in the model, the *individual impacts* of the environmental effects on the structural response, $\boldsymbol{\delta}_t^{\text{B|E}}$, are transformed in the PC space instead of

the environmental effects components. This enables to conserve the possibility of having different regression coefficients for the components of the model, a concept that was presented in Section 2.4. The vector of the individual impact of each environmental effect E_i is defined as

$$\boldsymbol{\delta}_t^{\text{B|E}} = \begin{bmatrix} \delta_t^{\text{B|E}_1} \\ \delta_t^{\text{B|E}_2} \\ \vdots \\ \delta_t^{\text{B|E}_i} \\ \vdots \\ \delta_t^{\text{B|E}_n} \end{bmatrix},$$

and it is obtained using the equation

$$\boldsymbol{\delta}_t^{\text{B|E}} = \mathbf{C}^{c,\text{B|E}} \mathbf{x}_t^{\text{E}}, \quad (3.8)$$

where

$$\mathbf{C}^{c,\text{B|E}} = \text{block diag} \left(\mathbf{C}^{c,\text{B|E}_1}, \mathbf{C}^{c,\text{B|E}_2}, \dots, \mathbf{C}^{c,\text{B|E}_i}, \dots, \mathbf{C}^{c,\text{B|E}_n} \right).$$

According to the proposed method, the transformed impacts are then

$$\begin{aligned} \boldsymbol{\delta}_t^{\text{B|E}^*} &= \boldsymbol{\phi}^{\text{PC}} \mathbf{P}^\top \boldsymbol{\delta}_t^{\text{B|E}} \\ &= \boldsymbol{\phi}^{\text{PC}} \mathbf{P}^\top \left(\mathbf{C}^{c,\text{B|E}} \mathbf{x}_t^{\text{E}} \right), \end{aligned} \quad (3.9)$$

where

$$\boldsymbol{\phi}^{\text{PC}} = \text{block diag} \left(\phi_1^{\text{PC}}, \phi_2^{\text{PC}}, \dots, \phi_p^{\text{PC}}, \dots, \phi_n^{\text{PC}} \right).$$

The equation for the transformed impacts using the proposed method with $m \leq n$ principal components is

$$\boldsymbol{\delta}_t^{\text{B|E}^*} = \begin{bmatrix} \delta_{1,t}^{\text{B|E}^*} \\ \delta_{2,t}^{\text{B|E}^*} \\ \vdots \\ \delta_{p,t}^{\text{B|E}^*} \\ \vdots \\ \delta_{m,t}^{\text{B|E}^*} \end{bmatrix} = \begin{bmatrix} \phi_1^{\text{PC}} P_{1,1} & \phi_1^{\text{PC}} P_{2,1} & \cdots & \phi_1^{\text{PC}} P_{i,1} & \cdots & \phi_1^{\text{PC}} P_{n,1} \\ \phi_2^{\text{PC}} P_{1,2} & \phi_2^{\text{PC}} P_{2,2} & \cdots & \phi_2^{\text{PC}} P_{i,2} & \cdots & \phi_2^{\text{PC}} P_{n,2} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \phi_p^{\text{PC}} P_{1,p} & \phi_p^{\text{PC}} P_{2,p} & \cdots & \phi_p^{\text{PC}} P_{i,p} & \cdots & \phi_p^{\text{PC}} P_{n,p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \phi_m^{\text{PC}} P_{1,m} & \phi_m^{\text{PC}} P_{2,m} & \cdots & \phi_m^{\text{PC}} P_{i,m} & \cdots & \phi_m^{\text{PC}} P_{n,m} \end{bmatrix} \begin{bmatrix} \delta_t^{\text{B|E}_1} \\ \delta_t^{\text{B|E}_2} \\ \vdots \\ \delta_t^{\text{B|E}_i} \\ \vdots \\ \delta_t^{\text{B|E}_n} \end{bmatrix}. \quad (3.10)$$

Note that $p \in \{1, 2, \dots, m\}$ defines the number on principal components in the model. The

analytical formulation of each individual transformed impact $\delta_{p,t}^{\mathbf{B}|\mathbf{E}^*}$ is

$$\delta_{p,t}^{\mathbf{B}|\mathbf{E}^*} = \overbrace{\phi_p^{\text{PC}} P_{1,p} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_1}}^{\text{contribution of } \mathbf{E}_1 \text{ to PC } p} + \overbrace{\phi_p^{\text{PC}} P_{2,p} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_2}}^{\text{contribution of } \mathbf{E}_2 \text{ to PC } p} + \cdots + \overbrace{\phi_p^{\text{PC}} P_{n,p} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_n}}^{\text{contribution of } \mathbf{E}_n \text{ to PC } p}. \quad (3.11)$$

However, because the objective is to integrate the PCA transformation in the existing formulation of BDLM, the individual transformed impacts in the PC space are broken down and the terms are reassembled by environmental effect observation. The transformed impacts of environmental effect observation \mathbf{E}_i in the original space is then a sum of its contribution to the m principal components

$$\begin{aligned} [\delta_t^{\mathbf{B}|\mathbf{E}_i}]_{\text{O}} &= \overbrace{\phi_1^{\text{PC}} P_{i,1} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_i}}^{\text{contribution of } \mathbf{E}_i \text{ to PC } 1} + \overbrace{\phi_2^{\text{PC}} P_{i,2} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_i}}^{\text{contribution of } \mathbf{E}_i \text{ to PC } 2} + \cdots + \overbrace{\phi_m^{\text{PC}} P_{i,m} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_i}}^{\text{contribution of } \mathbf{E}_i \text{ to PC } m} \\ &= (\phi_1^{\text{PC}} P_{i,1} + \phi_2^{\text{PC}} P_{i,2} + \cdots + \phi_m^{\text{PC}} P_{i,m}) \delta_t^{\mathbf{B}|\mathbf{E}_i} \\ &= \sum_{p=1}^m \phi_p^{\text{PC}} P_{i,p} \cdot \delta_t^{\mathbf{B}|\mathbf{E}_i} \\ &= \sum_{p=1}^m \phi_p^{\text{PC}} P_{i,p} (\mathbf{C}^{c,\mathbf{B}|\mathbf{E}_i} \mathbf{x}_t^{\mathbf{E}_i}), \end{aligned} \quad (3.12)$$

where $[\]_{\text{O}}$ indicates that the impact is in the original space. Then, the transformed impact of an environmental effect observation in the original space can be inserted in the BDLM formulation by modifying the observation matrix $\mathbf{C}^{c,\mathbf{B}|\mathbf{E}_i}$ by $\mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_i}$

$$\mathbf{C}^{c,\mathbf{B}|\mathbf{E}_i} \leftarrow \mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_i} = \mathbf{C}^{c,\mathbf{B}|\mathbf{E}_i} \cdot \sum_{p=1}^m \phi_p^{\text{PC}} \cdot P_{i,p} \quad 1 \leq m \leq n, \quad (3.13)$$

which quantifies the transformed impact of an environmental effect \mathbf{E}_i on the structural response \mathbf{B} . Note that there is a different coefficient for each environmental effect \mathbf{E}_i . Using these coefficients in the observation matrix \mathbf{C} of the dependent observation model is equivalent to using the transformation of Equation 3.9.

Figure 3.2 presents the flowchart of the method. First, the environmental effects need to be preprocessed before including them in the structural response model. To begin with, a model is built for each environmental effect observation. Each individual model is calibrated using MLE algorithm, and the hidden state variables are estimated using the Kalman filter and smoother. PCA's data matrix \mathbf{Z}^μ is built using Equation 3.6. Through PCA decomposition on the matrix \mathbf{Z}^μ , the PCA coefficient matrix \mathbf{P} and the explained vector $\boldsymbol{\varepsilon}$ are obtained. Equation 3.13 is used to include the transformed impact of the environmental effects models in the structural response model. The number of PCs included in the model can be reduced

by changing the value of m in Equation 3.13, which modifies the observation matrix \mathbf{C} for the structural response model, or equivalently by forcing the regression factors ϕ_p^{PC} to be equal to zero. The estimation of the structural response model's unknown parameters enables to obtain a set of optimized parameters, and so an optimized structural response model for the entire dataset, and to proceed to prediction. The portion of the flowchart within the dashed box represents the addition of the master thesis to the existing method, allowing to model a structural response with correlated environmental effects observations.

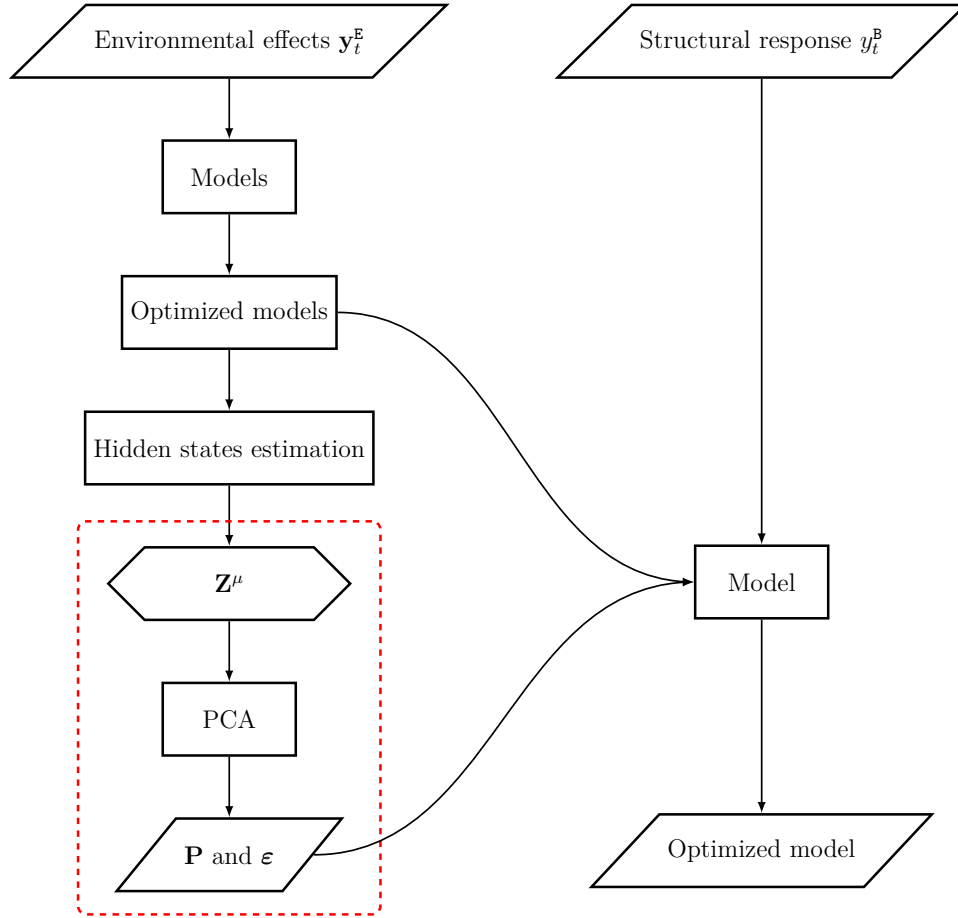


Figure 3.2 Flowchart representing the main steps of the proposed method. The dashed box contains the elements for the method proposed in this master thesis

CHAPTER 4 CASE STUDY

The case study presents an application of the method proposed in this thesis on SHM data acquired on a bridge located in Canada. This example presents the potential of the new method and discusses the effect of using only the few most important principal components versus using all of them. The prediction performances of the PCA-based method are also compared to the current version of BDLM.

4.1 Data Description

The case-study analyzes displacement and temperature data collected on a hyperstatic highway bridge located in Canada. The bridge is a prestressed reinforced concrete, cast in place structure. More specifically, a displacement dataset and four temperature datasets are employed for this example. Figure 4.1(a) presents a diagram of an elevation view of the structure and Figure 4.1(b) presents a cross-section to illustrate the location of the displacement and temperature sensors on the structure. The instrumented span and the total bridge are respectively approximately 30,5 meters and 106 meters long. The system of acquisition consists of a displacement sensor labeled d , and 4 temperature sensors named T1, T2, T3, and T4. The displacements are a longitudinal measurement of concrete expansion along a crack, expressed in millimetres and the temperatures are measured at different locations in the concrete structure, in degrees Celsius. The measurement precision for the displacement and temperature sensors is unknown. Both displacement and temperature data are recorded from November 2012 to October 2015 with a uniform time step of 1 hour, for a total of 25609 data points for each dataset. Note that there are some missing data in the dataset.

Figure 4.2(a) presents a superposition of the 4 raw temperature time series. They display a seasonal pattern in which the temperatures reach their maximum during summer and their minimum during winter. It is observed that the 4 temperature datasets are correlated to each other for the long-term scale, and a correlation is also observable for the short-term scale in the right section. However, some difference in the datasets can be observed on the short-term scale, as opposed to the long-term scale, where the 4 datasets are almost not differentiable. This can be explained by the position of the temperature sensor. Sensor T1 is the closest sensor to the upper surface, and so it is more sensitive to variations in air temperature and solar radiation, as opposed to sensors located deeper in the concrete. Figure 4.2(b) presents the raw displacement time series. Similarly to the 4 temperature time series, a yearly periodic pattern is apparent in the displacement time series : the displacements are

maximal during winter and minimal during summer. This behaviour of the displacement i.e., structural response, is explained by its dependency on the temperature i.e., environmental effect. There is a change in the amplitude of short-term variation depending on the time of the year, which is visible on both temperature and displacement data. The amplitude of short-term variation is larger during winter for all time series.

4.2 Description of the cases

For all the different cases, the training period is the first 2 years of data, or 17220 data points. The last year of data is used as the test dataset for validation. Different cases are studied, and they are defined in Table 4.1.

The first 4 cases employ 4 temperature observations with a different number of PCs to model the displacement observation d . The 5th to 8th cases are composed of the displacement observation modelled with a dependency on a single temperature observation at the time. The metric for comparing cases is the log-likelihood for the displacement observation only. As implied in Equation 2.19, the global log-likelihood for a model is the sum of the log-likelihood of

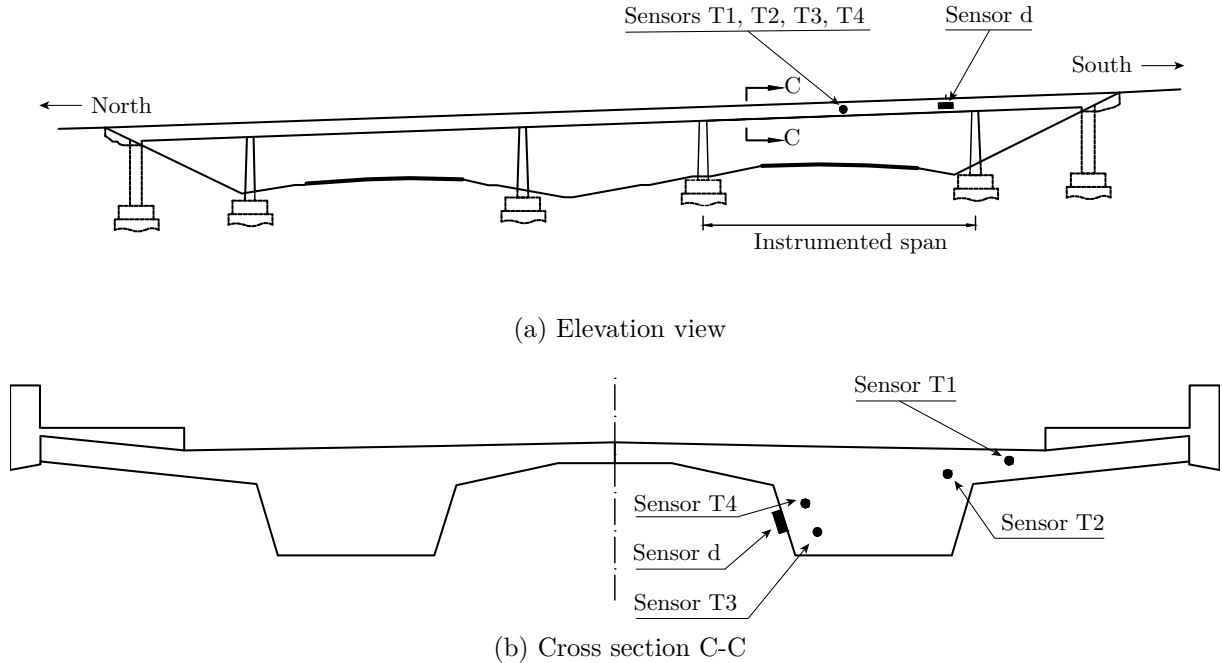


Figure 4.1 Diagrams showing the location of displacement sensor d and temperature sensors $T1$, $T2$, $T3$, and $T4$ on the bridge through (a) an elevation view and (b) a cross-section. A rectangle represents a displacement sensor and a circle represents a temperature sensors

its parts i.e., the displacement observation model and the temperature observation model(s). Because the cases have different combinations of temperature observation models, the impact of those temperature observations needs to be removed from the log-likelihood in order to have a consistent metric when comparing of the cases. Therefore, the cases are compared based on the log-likelihood of the optimized displacement observation model for the test period. Note that, for all cases, the model error is non-zero only for the autoregressive component σ_w^{AR} for simplification purposes and to reduce the number of unknown parameters

$$\sigma^{\text{LL}} = \sigma^{\text{P},S1} = \sigma^{\text{P},S2} = 0.$$

4.3 Models for cases 1 to 4 using PCA

This section presents the steps to obtain the optimized displacement models for cases 1 to 4. The process begins with the construction of the environmental effects observations models. Then, the PCA step is detailed, and finally the displacement model is build.

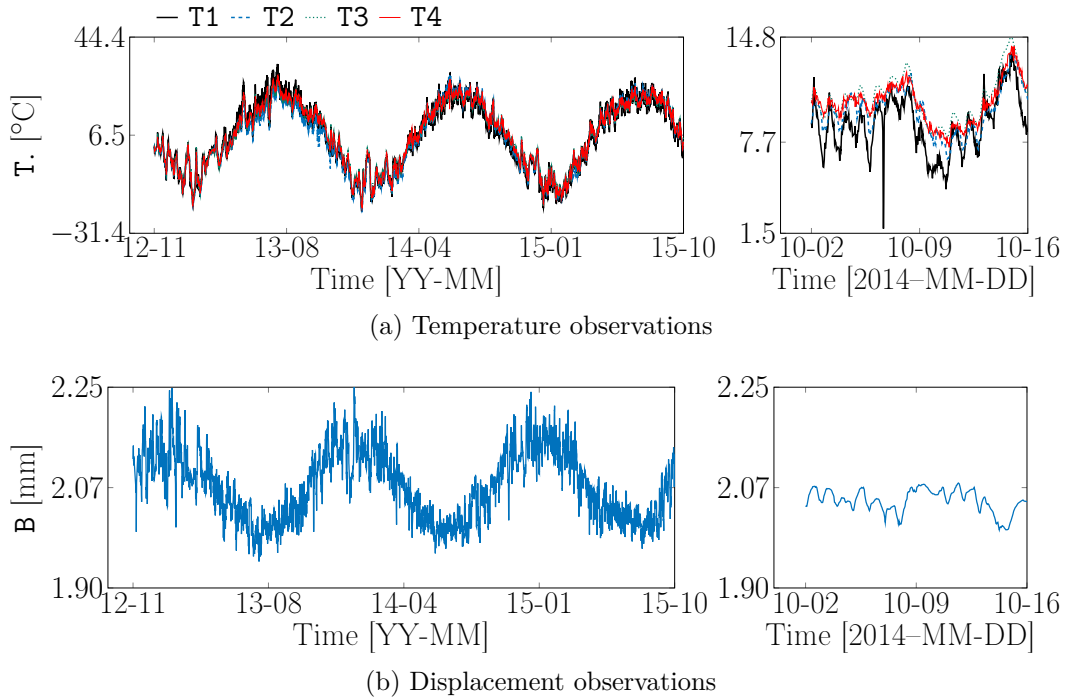


Figure 4.2 Raw data for the case study. The left section presents the entire dataset and the rights presents a 2-week period of data

Table 4.1 Table indicating the observations included in each case. Displacement observation is denoted d and temperature observations are T1, T2, T3, and T4. m refers to the number of principal components in the model

Case ID	d	T1	T2	T3	T4	m
1	✓	✓	✓	✓	✓	1
2	✓	✓	✓	✓	✓	2
3	✓	✓	✓	✓	✓	3
4	✓	✓	✓	✓	✓	4
5	✓	✓				-
6	✓		✓			-
7	✓			✓		-
8	✓				✓	-

4.3.1 Models for Environmental Effect Observations

For cases 1 to 4, the displacement model contains a displacement observation and 4 temperature observations. A model for an environmental effect observation refers to a model for only one given environmental effect observation and it constitutes in itself a portion of the displacement model.

Model Construction

In the raw temperature data presented in Figure 4.2(a), a similar behaviour is observed for the 4 temperature datasets so that a model with the same vector of hidden state variables is built for the 4 temperature observations. For a given environmental effect observation E_i , the vector of hidden state variables is composed of a local level (LL) to model the average temperature, a superposition of 3 periodic components (P) with periods of $P1 = 1$ day, $P2 = 365.24$ days to model the daily and seasonal fluctuations respectively, and $P3 = 1.0027$ days, and an autoregressive component (AR) to capture the time varying model errors. Note that the periodic component $P3$ is added to model the seasonal variation of amplitude in the daily periodic component, which is visible on Figure 4.2(a), by using the beat phenomenon (Roberts, 2016). The equation to calculate the period for the beat $P3$ is

$$P3 = \frac{1}{\frac{1}{P1} - \frac{1}{P2}} = \frac{1}{\frac{1}{1} - \frac{1}{365.24}} = 1.0027 \text{ days.} \quad (4.1)$$

Because the angular frequencies are needed in the model matrices, they are calculated

$$\begin{aligned}\omega^{\text{P1}} &= \frac{2\pi}{1} \\ \omega^{\text{P2}} &= \frac{2\pi}{365.24} \\ \omega^{\text{P3}} &= \frac{2\pi}{1.0027}.\end{aligned}\tag{4.2}$$

For the environmental effect observation models, the vector of hidden state variables is defined as

$$\mathbf{x}_t^{\text{E}_i} = \left[\overbrace{x_t^{\text{LL}}}^{\text{Local level}}, \overbrace{x_t^{\text{P1,S1}}, x_t^{\text{P1,S2}}}^{\text{Periodic P1}}, \overbrace{x_t^{\text{P2,S1}}, x_t^{\text{P2,S2}}}^{\text{Periodic P2}}, \overbrace{x_t^{\text{P3,S1}}, x_t^{\text{P3,S2}}}^{\text{Periodic P3}}, \overbrace{x_t^{\text{AR}}}^{\text{Autoregressive}} \right]^{\text{T}}.\tag{4.3}$$

The transition matrix is

$$\mathbf{A}^{\text{E}_i} = \text{block diag} \left(\mathbf{A}^{\text{LL}}, \mathbf{A}^{\text{P1}}, \mathbf{A}^{\text{P2}}, \mathbf{A}^{\text{P3}}, \mathbf{A}^{\text{AR}} \right),\tag{4.4}$$

where

$$\begin{aligned}\mathbf{A}^{\text{LL}} &= [1] \\ \mathbf{A}^{\text{P}} &= \begin{bmatrix} \cos(\omega^{\text{P}} \cdot \Delta t) & \sin(\omega^{\text{P}} \cdot \Delta t) \\ -\sin(\omega^{\text{P}} \cdot \Delta t) & \cos(\omega^{\text{P}} \cdot \Delta t) \end{bmatrix} \\ \mathbf{A}^{\text{AR}} &= [\phi^{\text{AR}}] \quad 0 \leq \phi^{\text{AR}} \leq 1.\end{aligned}$$

The observation matrix is

$$\begin{aligned}\mathbf{C}^{\text{E}_i} &= \begin{bmatrix} \mathbf{C}^{\text{LL}} & \mathbf{C}^{\text{P1}} & \mathbf{C}^{\text{P2}} & \mathbf{C}^{\text{P3}} & \mathbf{C}^{\text{AR}} \end{bmatrix} \\ \mathbf{C}^{\text{LL}} &= [1] \\ \mathbf{C}^{\text{P}} &= \begin{bmatrix} 1 & 0 \end{bmatrix} \\ \mathbf{C}^{\text{AR}} &= [1].\end{aligned}\tag{4.5}$$

Therefore,

$$\mathbf{C}^{\text{E}_i} = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix}.\tag{4.6}$$

The measurement error is distributed as

$$v^{\text{E}_i} \sim \mathcal{N}(0, \mathbf{R}^{\text{E}_i}),\tag{4.7}$$

where $\mathbf{R}^{E_i} = [(\sigma_v)^2]$ is the measurement variance. The model error is distributed as

$$\mathbf{w}^{E_i} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{E_i}), \quad (4.8)$$

where

$$\mathbf{Q}^{E_i} = \text{block diag} \left((\sigma^{\text{LL}})^2, \begin{bmatrix} (\sigma^{\text{P1},\text{S1}})^2 & 0 \\ 0 & (\sigma^{\text{P1},\text{S2}})^2 \end{bmatrix}, \begin{bmatrix} (\sigma^{\text{P2},\text{S1}})^2 & 0 \\ 0 & (\sigma^{\text{P2},\text{S2}})^2 \end{bmatrix}, \begin{bmatrix} (\sigma^{\text{P3},\text{S1}})^2 & 0 \\ 0 & (\sigma^{\text{P3},\text{S2}})^2 \end{bmatrix}, (\sigma^{\text{AR}})^2 \right). \quad (4.9)$$

For simplification purposes, it is assumed that only σ^{AR} is non-zero :

$$\mathbf{Q}^{E_i} = \text{block diag} \left(0, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, (\sigma^{\text{AR}})^2 \right). \quad (4.10)$$

Parameter and Hidden State Variable Estimation

For a given environmental effect E_i , the set of unknown parameters from the model matrices is

$$\mathcal{P}_{E_i} = \{\phi^{\text{AR}}, \sigma^{\text{AR}}, \sigma_v\}, \quad (4.11)$$

where ϕ^{AR} is the autocorrelation coefficient, σ^{AR} is the autocorrelation standard deviation, and σ_v is the observational error standard deviation. Parameter ϕ^{AR} varies from 0 to 1 and parameters σ^{AR} and σ_v are real numbers \mathbb{R}^+ . The MLE method presented in Section 2.6 is used to find the set of optimized parameters $\mathcal{P}_{E_i}^*$ for environmental effect E_i . The Kalman filter and smoother are used to estimate the hidden state variables distributions at every timestamp. The time dependent marginal distributions of the hidden state variables for the environmental effect observation $E_3 = \mathbf{T3}$ are presented in Appendix A. The models for the three other temperature observations are similar. The autoregressive components follow a stationary process for the four temperature models. The model predictions are presented with the raw data, and similar results are obtained for all the temperature observations.

The log-likelihood values over the test set for each temperature model are presented in Table 4.2. The log-likelihood can be used to compare models only if the models are based on the same dataset, for example to compare the performance of two models with a different combination of hidden states. On the other hand, $\mathbf{T3}$ might not be the observation that is the most correlated with the displacement observation \mathbf{d} because the temperature model performance is not related to the displacement model performance. Therefore, an observation different than $\mathbf{T3}$ could explain the displacement variations with more accuracy. Moreover, the likelihoods of the temperature observations have different orders of magnitude, which supports the point that the temperature observation log-likelihoods need to be removed in

order to compare the cases based on a comparable metric.

Table 4.2 Log-likelihood of the temperature models over the test set

Observation	Log-likelihood
T1	-5712
T2	765
T3	2727
T4	-2632

Principal Component Analysis

Prior to performing PCA on the environmental effects, the PCA's data matrix \mathbf{Z}^μ needs to be built using Equation 3.6. The terms in the equation are defined as

$$\begin{aligned} \mathbf{C}^{\mathbf{E}_i^*} &= \begin{bmatrix} \underbrace{0}_{\text{Local level}}, & \underbrace{1, 0}_{\text{Periodic P1}}, & \underbrace{1, 0}_{\text{Periodic P2}}, & \underbrace{1, 0}_{\text{Periodic P3}}, & \underbrace{1}_{\text{Autoregressive}} \end{bmatrix} \\ \boldsymbol{\mu}_{t|T}^{\mathbf{E}_i} &= \begin{bmatrix} \underbrace{\mu_{t|T}^{\text{LL}}}_{\text{Local level}}, & \underbrace{\mu_{t|T}^{\text{P1,S1}}, \mu_{t|T}^{\text{P1,S2}}}_{\text{Periodic P1}}, & \underbrace{\mu_{t|T}^{\text{P2,S1}}, \mu_{t|T}^{\text{P2,S2}}}_{\text{Periodic P2}}, & \underbrace{\mu_{t|T}^{\text{P3,S1}}, \mu_{t|T}^{\text{P3,S2}}}_{\text{Periodic P3}}, & \underbrace{\mu_{t|T}^{\text{AR}}}_{\text{Autoregressive}} \end{bmatrix}^\top. \end{aligned} \quad (4.12)$$

In this example, because the four temperature models have the same components, the formulation is the same for all the environmental effect observations i.e.,

$$z_{t,\mathbf{E}_i}^\mu = \mathbf{C}^{\mathbf{E}_i^*} \boldsymbol{\mu}_{t|T}^{\mathbf{E}_i} = \mu_{t|T}^{\text{P1,S1}} + \mu_{t|T}^{\text{P2,S1}} + \mu_{t|T}^{\text{P3,S1}} + \mu_{t|T}^{\text{AR}}. \quad (4.13)$$

Figure 4.3 illustrates the construction of \mathbf{Z}^μ matrix with the specifications of this example.

The size of \mathbf{Z}^μ matrix is $[T \times n] = [25609 \times 4]$. PCA enables to compute the coefficient matrix \mathbf{P} , of dimensions $[n \times n] = [4 \times 4]$, and the explained vector $\boldsymbol{\varepsilon}$

$$\begin{aligned} \mathbf{P} &= \begin{bmatrix} 0.511 & 0.708 & -0.483 & 0.055 \\ 0.505 & -0.699 & -0.499 & -0.083 \\ 0.486 & 0.060 & 0.523 & -0.698 \\ 0.497 & -0.078 & 0.493 & 0.710 \end{bmatrix} \\ \boldsymbol{\varepsilon} &= [99.37 \quad 0.416 \quad 0.186 \quad 0.032]^\top. \end{aligned} \quad (4.14)$$

From the values in the explained vector $\boldsymbol{\varepsilon}$, the first principal component explains 99.37% of the four temperature observations total variance. This suggests that the temperature obser-

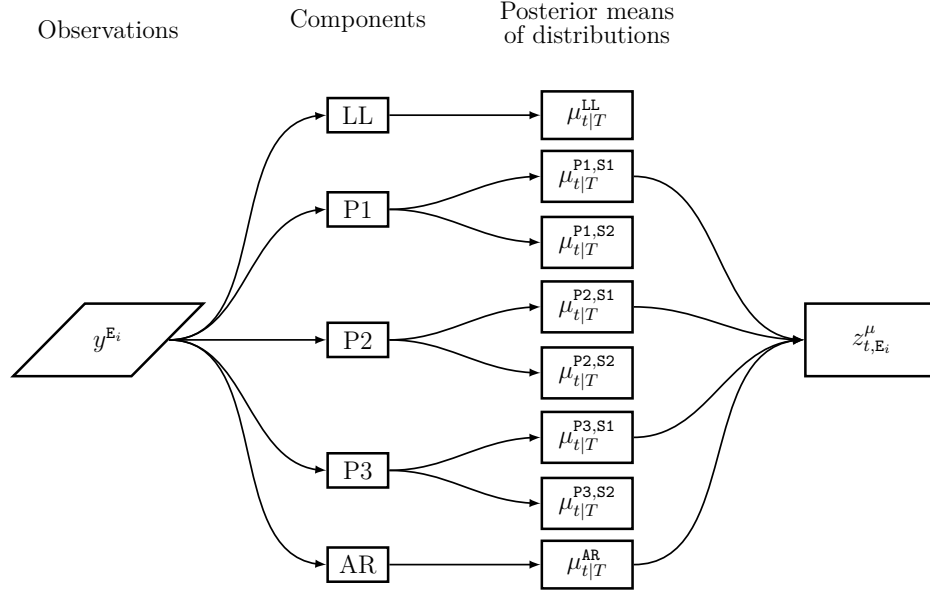


Figure 4.3 Diagram of the steps with an environmental effect observations to build \mathbf{Z}^μ matrix

variations are strongly correlated as one axis describes the majority of variance. Determining the percentage of the total variance explained by the first m PC is made possible by summing the m first values from the explained vector $\boldsymbol{\varepsilon}$. For this example, using 2 PCs contains 99.79% of the temperature observations variance and 3 PCs explains 99.97%. Using the total number of PCs always explains 100% of the variance. Removing some PCs in the model is a way to reduce the number of unknown parameters as each PC p has its own scaling factor ϕ_p^{PC} in the structural response model. For m PCs, there are m parameters, $\{\phi_1^{\text{PC}}, \phi_2^{\text{PC}}, \dots, \phi_m^{\text{PC}}\}$, which correspond to the scaling factors or regression coefficient of each principal component.

4.3.2 Model for Structural Response Observation

The model for the structural response observation corresponds to the model describing the behaviour of the displacement observation including the temperature dependencies. This section details the dependencies and how the environmental effect models are assembled with the displacement's own components to form the displacement model. The sections presents the different hypotheses that simplify the models, and details the model construction, and the parameter and hidden state estimation.

Hypotheses

To simplify the example, some hypotheses are made concerning the displacement model with the dependencies on temperature. First, the optimized parameters of the temperature observation models used to compute the PCA coefficient matrix are assumed to be of the same values in the displacement model. Second, the model matrices $\{\mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}\}$ are assumed to be constant over time. Third, it is assumed that the regression coefficients of a given type of component are equal for the 4 temperature observations. As an example, the regression coefficients for the autoregressive components are

$$\phi_{\text{AR}}^{\text{B|E}_1} = \phi_{\text{AR}}^{\text{B|E}_2} = \phi_{\text{AR}}^{\text{B|E}_3} = \phi_{\text{AR}}^{\text{B|E}_4},$$

and the regression coefficient is generalized and labelled $\phi_{\text{AR}}^{\text{B|E}}$. This enables to constrain 12 regression coefficients and to reduce the number of unknown parameters. Then, on Figure 4.2, one can note that the maximum displacements occur during the coldest periods and, inversely, the minimal displacements occur during the hottest period of the year. This relationship demonstrates a negative correlation between the displacement and temperature on the long term scale. On the opposite, on Figure 4.4, which shows a close-up of the raw data, a positive correlation between the displacement and temperature observations can be noted. For BDLMs, this suggests that the regression coefficients are positive for the short term components (daily periodic component and autoregressive component) and negative for the long term components (seasonal periodic component). To account for the phenomenon, the regression coefficients are set to be independent for all the components.

Model Construction

The observation vector for the displacement model with the environmental effect dependencies is defined as

$$\mathbf{y}_t = \begin{bmatrix} y_t^{\text{B}} \\ \mathbf{y}_t^{\text{E}} \end{bmatrix}, \quad (4.15)$$

where

$$y_t^{\text{B}} = \begin{bmatrix} y_t^{\text{d}} \end{bmatrix}$$

$$\mathbf{y}_t^{\text{E}} = \begin{bmatrix} y_t^{\text{T1}} & y_t^{\text{T2}} & y_t^{\text{T3}} & y_t^{\text{T4}} \end{bmatrix}^{\text{T}}.$$

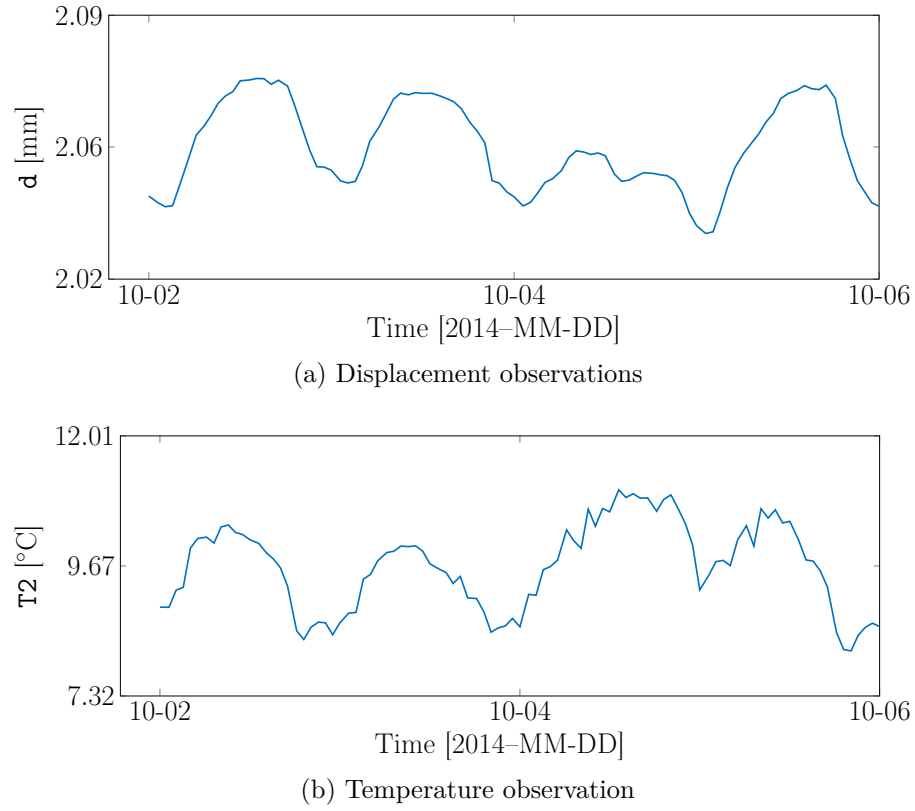


Figure 4.4 Displacement d and temperature $T2$ raw data showing a positive correlation between displacement and temperature observations on the short term scale

For the displacement observation model, there are $n = 4$ environmental effect observations of a same type. The dependence matrix is

$$\mathbf{D} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad (4.16)$$

which indicates that the first observation i.e., the displacement, is dependent on the other four temperature observations. The displacement observation has its own components, which are set to a local level (^{LL}) and an autoregressive component (^{AR}). The vector of hidden state variables for the displacement model is

$$\mathbf{x}_t = \begin{bmatrix} \mathbf{x}_t^{\text{B}} \\ \mathbf{x}_t^{\text{E}_1} \\ \mathbf{x}_t^{\text{E}_2} \\ \mathbf{x}_t^{\text{E}_3} \\ \mathbf{x}_t^{\text{E}_4} \end{bmatrix}, \quad (4.17)$$

where

$$\mathbf{x}_t^{\text{B}} = \left[\overbrace{x_t^{\text{LL}}}^{\text{Local level}}, \overbrace{x_t^{\text{AR}}}^{\text{Autoregressive}} \right]^{\text{T}}.$$

The transition matrix \mathbf{A} is defined as

$$\mathbf{A} = \text{block diag} \left(\mathbf{A}^{\text{B}}, \mathbf{A}^{\text{E}_1}, \mathbf{A}^{\text{E}_2}, \mathbf{A}^{\text{E}_3}, \mathbf{A}^{\text{E}_4} \right), \quad (4.18)$$

where

$$\mathbf{A}^{\text{B}} = \begin{bmatrix} \mathbf{A}^{\text{LL}} & 0 \\ 0 & \mathbf{A}^{\text{AR}} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & \phi^{\text{AR}} \end{bmatrix}.$$

Because of the dependency between the displacement and temperature observations, non-zero terms are present outside the diagonal in the observation matrix \mathbf{C} of the structural

response model. The observation matrix is

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}^{\mathbf{B}} & \mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_1} & \mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_2} & \mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_3} & \mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_4} \\ 0 & \mathbf{C}^{\mathbf{E}_1} & 0 & 0 & 0 \\ 0 & 0 & \mathbf{C}^{\mathbf{E}_2} & 0 & 0 \\ 0 & 0 & 0 & \mathbf{C}^{\mathbf{E}_3} & 0 \\ 0 & 0 & 0 & 0 & \mathbf{C}^{\mathbf{E}_4} \end{bmatrix}, \quad (4.19)$$

where

$$\begin{aligned} \mathbf{C}^{\mathbf{B}} &= \begin{bmatrix} \mathbf{C}^{\text{LL}} & \mathbf{C}^{\text{AR}} \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix} \\ \mathbf{C}_{\text{PC}}^{c,\mathbf{B}|\mathbf{E}_i} &= \begin{bmatrix} \phi_{\text{LL}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{P1,S1}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{P1,S2}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{P2,S1}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{P2,S2}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{P3,S1}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{P3,S2}}^{\mathbf{B}|\mathbf{E}} & \phi_{\text{AR}}^{\mathbf{B}|\mathbf{E}} \end{bmatrix} \cdot \sum_{p=1}^m \phi_p^{\text{PC}} \cdot P_{i,p} \\ &= \begin{bmatrix} 0 & \phi_{\text{P1,S1}}^{\mathbf{B}|\mathbf{E}} & 0 & \phi_{\text{P2,S1}}^{\mathbf{B}|\mathbf{E}} & 0 & \phi_{\text{P3,S1}}^{\mathbf{B}|\mathbf{E}} & 0 & \phi_{\text{AR}}^{\mathbf{B}|\mathbf{E}} \end{bmatrix} \cdot \sum_{p=1}^m \phi_p^{\text{PC}} \cdot P_{i,p}, \end{aligned}$$

using Equation 3.13. The observational errors are

$$\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}), \quad (4.20)$$

where

$$\begin{aligned} \mathbf{R} &= \text{block diag}(\mathbf{R}^{\mathbf{B}}, \mathbf{R}^{\mathbf{E}_1}, \mathbf{R}^{\mathbf{E}_2}, \mathbf{R}^{\mathbf{E}_3}, \mathbf{R}^{\mathbf{E}_4}) \\ \mathbf{R}^{\mathbf{B}} &= [(\sigma_v^{\mathbf{B}})^2] \\ \mathbf{v} &= \begin{bmatrix} v^{\mathbf{B}} & v^{\mathbf{E}_1} & v^{\mathbf{E}_2} & v^{\mathbf{E}_3} & v^{\mathbf{E}_4} \end{bmatrix}^{\top}. \end{aligned}$$

The model errors are

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}), \quad (4.21)$$

where

$$\begin{aligned} \mathbf{Q} &= \text{block diag}(\mathbf{Q}^{\mathbf{B}}, \mathbf{Q}^{\mathbf{E}_1}, \mathbf{Q}^{\mathbf{E}_2}, \mathbf{Q}^{\mathbf{E}_3}, \mathbf{Q}^{\mathbf{E}_4}) \\ \mathbf{Q}^{\mathbf{B}} &= \begin{bmatrix} (\sigma^{\text{LL}})^2 & 0 \\ 0 & (\sigma^{\text{AR}})^2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & (\sigma^{\text{AR}})^2 \end{bmatrix} \\ \mathbf{w} &= \begin{bmatrix} \mathbf{w}^{\mathbf{B}} & \mathbf{w}^{\mathbf{E}_1} & \mathbf{w}^{\mathbf{E}_2} & \mathbf{w}^{\mathbf{E}_3} & \mathbf{w}^{\mathbf{E}_4} \end{bmatrix}^{\top}. \end{aligned}$$

because only σ^{AR} is non-zero.

Parameters and Hidden State Variable Estimation

At this point, the unknown parameters are

$$\mathcal{P} = \{ \overbrace{\phi^{\text{AR}}, \sigma^{\text{AR}}, \sigma_v^{\text{B}}}^{\text{Parameters for B}}, \overbrace{\phi_{\text{AR}}^{\text{B|E}}, \phi_{\text{P1,S1}}^{\text{B|E}}, \phi_{\text{P2,S1}}^{\text{B|E}}, \phi_{\text{P3,S1}}^{\text{B|E}}}^{\text{Regression coefficients}}, \overbrace{\phi_1^{\text{PC}}, \phi_2^{\text{PC}}, \dots, \phi_m^{\text{PC}}}^{\text{PC scaling factors}} \}, \quad (4.22)$$

where ϕ^{AR} is the autocorrelation coefficient varying from 0 to 1, and σ^{AR} and σ_v^{B} are the autocorrelation and observational error standard deviation respectively. Using MLE, the optimized values of the parameters \mathcal{P}^* are estimated based on the data from the training period. The Kalman filter and smoother enable to estimate the hidden state variables distributions for the entire set of hidden state variables. For case 2, the distributions of the two components of the displacement observation (local level and autoregressive component) are presented Appendix B. The model for the displacement is also illustrated in Appendix B. The models for cases 1, 3 and 4 are similar. In the autoregressive components of the models for cases 1 to 4, there is no clear pattern remaining in the signal and it follows stationary process, which indicates that the model contains all the apparent components.

4.4 Models for cases 5 to 8 not using PCA

This section presents the models for cases 5 to 8. For those cases, the model contains a displacement observation and one temperature observation. The temperature observation models are defined in Section 4.3.1.

4.4.1 Model for Structural Response Observation

Model Construction

For the 5th to 8th cases, where PCA is not used, the observation vector for the displacement model using a given environmental effect observation \mathbf{E}_i is

$$\mathbf{y}_t = \begin{bmatrix} y_t^{\text{B}} \\ y_t^{\text{E}} \end{bmatrix} = \begin{bmatrix} y_t^{\text{d}} \\ y_t^{\text{Ti}} \end{bmatrix}. \quad (4.23)$$

The dependence matrix is

$$\mathbf{D} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \quad (4.24)$$

which indicates that the temperature observation has an impact on the displacement observation. For the displacement observation, the components are defined in Section 4.3.2. The

vector of hidden state variables is

$$\mathbf{x}_t = \left[\overbrace{x_t^{\text{LL}}, x_t^{\text{AR}}}^{\text{Displacement}}, \overbrace{x_t^{\text{LL}}, x_t^{\text{P1,S1}}, x_t^{\text{P1,S2}}, x_t^{\text{P2,S1}}, x_t^{\text{P2,S2}}, x_t^{\text{P3,S1}}, x_t^{\text{P3,S2}}, x_t^{\text{AR}}}^{\text{Temperature}} \right]^\top. \quad (4.25)$$

The model matrices are

$$\begin{aligned} \mathbf{A} &= \text{block diag} \left(\mathbf{A}^{\text{B}}, \mathbf{A}^{\text{E}_i} \right) \\ \mathbf{C} &= \begin{bmatrix} \mathbf{C}^{\text{B}} & \mathbf{C}^{\text{B}|\text{E}_i} \\ 0 & \mathbf{C}^{\text{E}_i} \end{bmatrix}, \end{aligned} \quad (4.26)$$

where

$$\begin{aligned} \mathbf{C}^{\text{B}|\text{E}_i} &= \begin{bmatrix} \phi_{\text{LL}}^{\text{B}|\text{E}_i} & \phi_{\text{P1,S1}}^{\text{B}|\text{E}_i} & \phi_{\text{P1,S2}}^{\text{B}|\text{E}_i} & \phi_{\text{P2,S1}}^{\text{B}|\text{E}_i} & \phi_{\text{P2,S2}}^{\text{B}|\text{E}_i} & \phi_{\text{P3,S1}}^{\text{B}|\text{E}_i} & \phi_{\text{P3,S2}}^{\text{B}|\text{E}_i} & \phi_{\text{AR}}^{\text{B}|\text{E}_i} \end{bmatrix} \\ &= \begin{bmatrix} 0 & \phi_{\text{P1,S1}}^{\text{B}|\text{E}_i} & 0 & \phi_{\text{P2,S1}}^{\text{B}|\text{E}_i} & 0 & \phi_{\text{P3,S1}}^{\text{B}|\text{E}_i} & 0 & \phi_{\text{AR}}^{\text{B}|\text{E}_i} \end{bmatrix}. \end{aligned}$$

The measurement and model errors are

$$\begin{aligned} \mathbf{v} &\sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \\ \mathbf{w} &\sim \mathcal{N}(\mathbf{0}, \mathbf{Q}), \end{aligned} \quad (4.27)$$

where

$$\begin{aligned} \mathbf{v} &= \begin{bmatrix} v^{\text{B}} \\ v^{\text{E}_i} \end{bmatrix} \\ \mathbf{R} &= \text{block diag} \left(\mathbf{R}^{\text{B}}, \mathbf{R}^{\text{E}_i} \right) \\ \mathbf{w} &= \begin{bmatrix} \mathbf{w}^{\text{B}} \\ \mathbf{w}^{\text{E}_i} \end{bmatrix} \\ \mathbf{Q} &= \text{block diag} \left(\mathbf{Q}^{\text{B}}, \mathbf{Q}^{\text{E}_i} \right). \end{aligned}$$

Parameter and Hidden State Variable Estimation

For a given case with a given environmental effect observation E_i , the unknown parameters from the model matrices are

$$\mathcal{P} = \left\{ \overbrace{\phi_{\text{AR}}^{\text{B}}, \sigma_{\text{AR}}^{\text{B}}, \sigma_v^{\text{B}}}^{\text{Parameters for B}}, \overbrace{\phi_{\text{P1,S1}}^{\text{B}|\text{E}_i}, \phi_{\text{P2,S1}}^{\text{B}|\text{E}_i}, \phi_{\text{P3,S1}}^{\text{B}|\text{E}_i}, \phi_{\text{AR}}^{\text{B}|\text{E}_i}}^{\text{Regression coefficients}} \right\}. \quad (4.28)$$

Using MLE method and the first two years of data as training period, the optimized parameters \mathcal{P}^* are estimated. The Kalman filter and smoother enable to estimate the hidden state variables.

4.5 Results

All the cases generated estimations for the displacement that are similar and in good agreement with the data. Their difference is easier to quantify using the log-likelihoods, calculated using Equation 2.20. For cases 1 to 8, the log-likelihood of the displacement observation during the test period are presented in Table 4.3. All the cases have a log-likelihood in the same order of magnitude. For cases 1 to 4, the highest log-likelihood is obtained by case 2, followed by cases 3 and 4. The lowest log-likelihood is obtained by case 1. The log-likelihoods of cases 5 to 8 are between the log-likelihood of cases 1 and 4, which are the cases using PCA with the two lowest log-likelihoods. For cases 5 to 8, the case with the highest log-likelihood is case 8, and the lowest is case 7.

Table 4.3 Displacement observation log-likelihood for the studied cases using PCA

Case ID	m	% variance	Temperature observations	Log-likelihood for test period
1	1	99.37	T1 to T4	35 837
2	2	99.78	T1 to T4	36 974
3	3	99.97	T1 to T4	36 967
4	4	100	T1 to T4	36 951
5	-	100	T1	36 591
6	-	100	T2	36 692
7	-	100	T3	35 848
8	-	100	T4	36 935

4.6 Discussion

When comparing the log-likelihood of the displacement observation model for the test period from Table 4.3, it can be noted that the log-likelihood significantly increases from 1 to 2 PCs, but decreases when adding more PCs. In other words, adding the second PC increases the prediction capacity of the displacement observation model, but adding the last 2 PCs, which jointly explain a total of 0.218% of the total variance, does not. This suggests two things: (1) the second PC carries information that is useful to predict the displacement and (2) the last 2 PCs of the temperature observations describe processes that are not useful to predict the displacement. Because a PC does not necessarily have a physical meaning, it is hard to interpret the information carried by this PC. However, the second PC does carry information that is not available when only 1 sensor is used. This suggests that the first PC carries the information that is shared by all the temperature observations. As presented in Figure

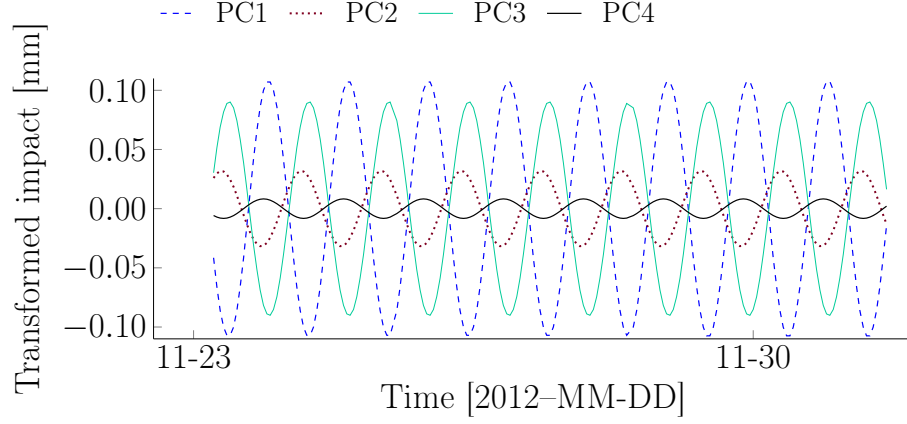
4.5, which represents the time dependent transformed impacts of the various components in the PC space for case 4, the four temperature sensors are strongly correlated based on the annual cycle and this correlation is mostly represented by PC 1, because the only PC with significant amplitude is PC 1 in Figure 4.5(b). PCs 2, 3 and 4 have negligible information concerning the seasonal cycles. On the other hand, the second PC contains the information that differentiates the temperature observations from each other on the short term basis because this information is less correlated. This last phenomenon is visible in Figure 4.5(a), (c), and (d), where it can be observed that information from the impacts of the autoregressive component, the daily periodic component and the beat periodic component is contained in each PC. That information is, for example, the phase shift between sensors which enables to take into account the spatial variability. The strong correlation on the yearly cycle and the weaker correlation in the short-term basis are visible on the raw data on Figure 4.2. Note that PC 1 that explains a large variance does not necessarily carries the important information for predicting the displacement i.e., PC 1 might not be the PC that is the most correlated to the displacement. On the other hand, the 3rd and 4th PCs carry information that explains a small portion of the variance. This information is not helpful to predict the displacement, as adding them does not increases the log-likelihood. An hypothesis about this phenomenon is that PC 3 and 4 have either a negligible amplitude, as seen with the 4th PC in Figure 4.5(a), (b), (c) and (d), or the impact is opposite to the others PCs and the PCs cancel out, for instance, PC 1 and 3 in Figure 4.5(a) are out of phase. Therefore, for this case study, two PCs are sufficient to extract the useful information from the temperature observations.

Then, the model with one temperature observation that performs the best for describing the displacement observation is the model with temperature observation T4. However, the log-likelihood for case 8 is lower than when including two to four PCs, from Table 4.3. In other words, a model using the proposed method with more than one PC performs better than any model with only one temperature observation. This suggests that the proposed method, when using all the available data, leads to better results and prediction capacities for a structural response model than using any temperature observation alone. The number of PCs in the model also has an impact on the prediction capacity, as explained previously, and it should be optimized. In brief, the case study demonstrates that the method has the potential to include environmental effect datasets in a model and to increase the prediction capacity, while having the possibility to remove the useless information through the number of PCs included.

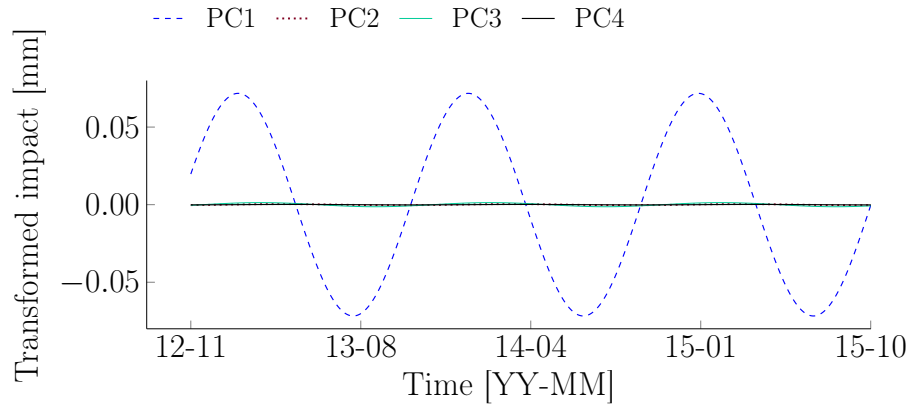
For cases 1 to 8, the optimized parameters are not presented because they are hardly comparable. In cases 5 to 8, the regression coefficients $\{\phi_{P1,S1}^{B|E_i}, \phi_{P2,S1}^{B|E_i}, \phi_{P3,S1}^{B|E_i}, \phi_{AR}^{B|E_i}\}$ are straightforward. However, for cases 1 to 4, several values of PC scaling factors can lead to similar results. In

addition to allowing to handle correlated external conditions, the proposed method supports missing data. If a data point is missing during the training period, the PCA can still be processed as the data matrix \mathbf{Z}^μ is build using the estimation of the hidden state variables from the Kalman filter and smoother, which is a prediction when no data is recorded for this timestamp. When a data point is missing once the training period is finished, the prediction from the Kalman filter without the update step is used. Therefore, the proposed method supports missing data. Using the data matrix \mathbf{Z}^μ also allow to handle outliers, as the estimation prevents from having aberrant data, and to remove the measurement error from the original data.

If different types of environmental effects are observed and included in a dataset, they could all be included in a model using the proposed method to allow increasing the prediction capacity. Some additional pre-processing is necessary. Examples of different types of en-

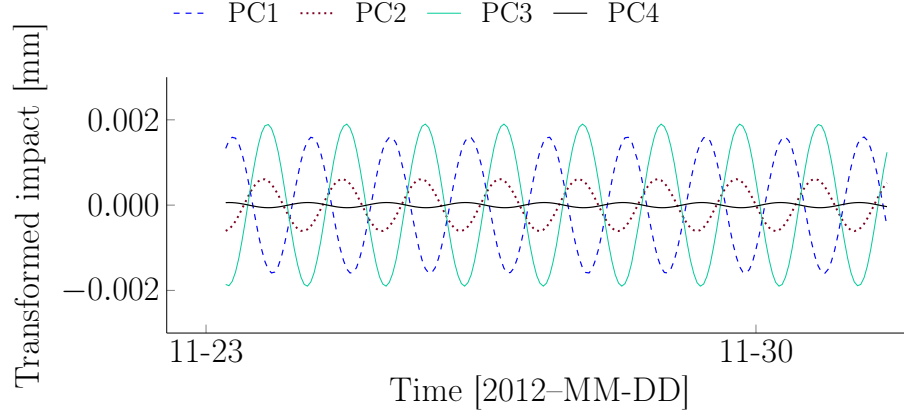


(a) First hidden state variable of daily periodic component $P1 = 1$ day

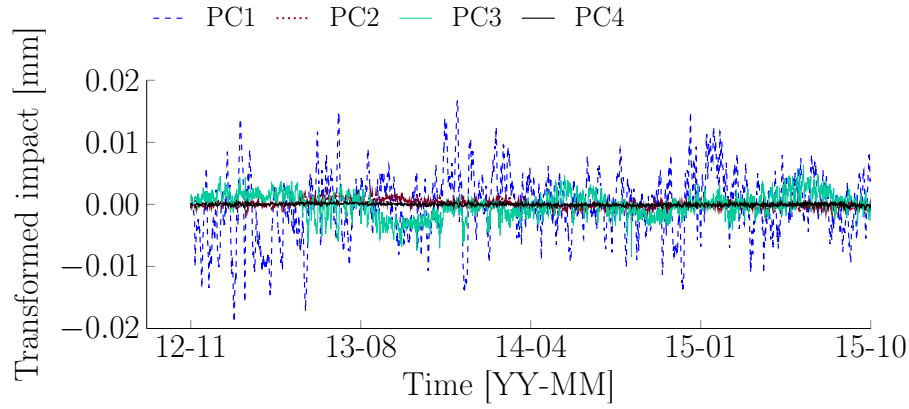


(b) First hidden state variable of daily periodic component $P2 = 365.24$ days

Figure 4.5 Transformed impact of the components from the 4 temperature observations in case 4



(c) First hidden state variable of daily periodic component $P3 = 1.0027$ days



(d) Autoregressive component

Figure 4.5 Transformed impact of the components from the 4 temperature observations in case 4

environmental effects are humidity, wind speed, and sun radiation, among others, and those other environmental effects could be correlated. The pre-processing is needed to treat all the environmental effects at once and to remove the scale effect. The pre-processing consists of dividing the dataset into sub-datasets where a given sub-dataset only has data from one type of environmental effect, and then normalizing the sub-datasets. Then the sub-datasets can be joined to form the data matrix \mathbf{Z}^μ , and PCA can be processed. This modification to the method only modifies the pre-processing step, as the normalization constant for each subset is transferred to the regression coefficients further on.

In brief, the case-study demonstrates that the proposed method is able to handle the dependency between a structural response and multiple correlated environmental effect observations.

CHAPTER 5 CONCLUSION AND RECOMMENDATIONS

5.1 Summary of work

This master thesis presents a solution to tackle the issue of indeterminate systems caused by correlated environmental effect observations using *Principal Component Analysis* and adapted to *Bayesian Dynamic Linear Models*. For the purpose of illustrating the method's capacity, an example is built using hourly data of displacement and four temperature observations from a highway concrete bridge located in Canada. Four models are built using all the temperature observations and 1 to 4 PCs, and four models are built using one temperature observation at the time. With this dataset, using 2 to 4 principal components with the PCA method produces models with better prediction capacities than the models without the proposed method. The proposed method also increases the robustness and reliability as redundant sensors are necessary, and a simultaneous failure of multiple sensors is unlikely.

Thus, in general, if one is confronted to a dataset with multiple environmental effect sensors, for example temperature, humidity, and solar radiation, the recommendation is to include all the available data in the model. The method enables the user to avoid discarding datasets that include information relevant for explaining the structural response.

5.2 Limitations

A first limitation is in the decision concerning the number of PCs to be included in a structural response model. The process to determine this number is not defined in this master thesis. For this master thesis, the methodology is to test all the possibility, which is time-consuming. The general recommendation from this thesis is to select more than 1 PC but not the entire set of PCs. Hua et al. and Ni et al. also arrived to the conclusion that a model performs better without including all PCs (Hua et al., 2007; Ni et al., 2006). A second limitation concerns the implementation of the method. The described method needs manipulations from the user between environmental effect model optimization and structural response model optimization, and the method is therefore not automatized yet. However, the method could be modified to be easier to implement by adding some assumptions and simplifications. For example, the raw data could be used to form the PCA's data matrix \mathbf{Z} . In this case, the entire set of PC should not be included in the model to limit the presence of observational noise, and missing data has to be removed. A third limitation is the hypotheses concerning the regression coefficients. In the case study, the regression coefficient for a given component

is assumed to be equal for all the temperature observations. This simplifies the model and limits the number of unknown parameters, but it also limits the flexibility. However, the results still indicates that the method increases the prediction capacity, which suggests that this hypothesis is justified.

A last limitation concerns the calculation of the coefficient matrix \mathbf{P} . It is computed using PCA's data matrix \mathbf{Z}^μ , which is composed of the expected value of the hidden state variables having an impact on the displacement. Using the expected value and not the distribution is a simplification that is allowed if and only if the variance is small, as the variance should be taken into account otherwise. In the case-study presented in this thesis, the hidden state variable variance is small to negligible, and it makes the hypothesis adequate.

5.3 Future work

Further studies are needed to define precisely the methodology for selecting an appropriate number of principal components. Then, the theory concerning the use of different environmental effect types has to be tested with real data, as this kind of data was not available for this study. It would also be interesting to do a more extended research to explore the impacts of thermal inertia in similar cases, which causes a delay in the structural response to the external condition and temperature gradients in the structure.

REFERENCES

- H. Abdi and L. J. Williams, “Principal Component Analysis”, *Wiley interdisciplinary reviews: computational statistics*, vol. 2, no. 4, pp. 433–459, 2010.
- M. Enckell, B. Glisic, F. Myrvoll, and B. Bergstrand, “Evaluation of a Large-scale Bridge Strain, Temperature and Crack Monitoring with Distributed Fibre Optic Sensors”, *Journal of Civil Structural Health Monitoring*, vol. 1, no. 1-2, pp. 37–46, 2011.
- A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin, *Bayesian Data Analysis*, 3rd Ed. CRC Press, 2014.
- J.-A. Goulet, “Bayesian Dynamic Linear Models for Structural Health Monitoring”, *Structural Control and Health Monitoring*, vol. 24, pp. e2035–n/a, 2017. DOI: 10.1002/stc.2035
- J. A. Goulet and K. Koo, “Empirical Validation of Bayesian Dynamic Linear Models in the Context of Structural Health Monitoring”, *ASCE, Journal of Bridge Engineering*, vol. 23, no. 2, p. 05017017, 2018.
- X. Hua, Y. Ni, J. Ko, and K. Wong, “Modeling of Temperature–Frequency Correlation using Combined Principal Component Analysis and Support Vector Regression Technique”, *ASCE, Journal of Computing in Civil Engineering*, vol. 21, no. 2, pp. 122–135, 2007.
- I. Jolliffe, “Principal Components in Regression Analysis”, in *Principal Component Analysis*, 2nd Ed. Springer, 2002, pp. 167–198.
- P. Léger and M. Leclerc, “Hydrostatic, Temperature, Time-displacement Model for Concrete Dams”, *ASCE, Journal of Engineering Mechanics*, vol. 133, no. 3, pp. 267–277, 2007. DOI: 10.1061/(ASCE)0733-9399(2007)133:3(267)
- H.-N. Li, L. Ren, Z.-G. Jia, T.-H. Yi, and D.-S. Li, “State-of-the-art in Structural Health Monitoring of Large and Complex Civil Infrastructures”, *Journal of Civil Structural Health Monitoring*, vol. 6, no. 1, pp. 3–16, 2016.
- M. P. Limongelli, E. Chatzi, M. Döhler, G. Lombaert, and E. Reynders, “Towards Extraction of Vibration-based Damage Indicators”, in *EWSHM-8th European Workshop on Structural Health Monitoring*, 2016.

J. Lynch and K. Loh, “A Summary Review of Wireless Sensors and Sensor Networks for Structural Health Monitoring”, *Shock and Vibration Digest*, vol. 38, no. 2, pp. 91–130, 2006.

F. Magalhães, A. Cunha, and E. Caetano, “Vibration Based Structural Health Monitoring of an Arch Bridge: from Automated OMA to Damage Detection”, *Mechanical Systems and Signal Processing*, vol. 28, pp. 212–228, 2012.

M. Malekzadeh, G. Atia, and F. Catbas, “Performance-based Structural Health Monitoring through an Innovative Hybrid Data Interpretation Framework”, *Journal of Civil Structural Health Monitoring*, vol. 5, no. 3, pp. 287–305, 2015.

J. Mata, “Interpretation of Concrete Dam Behaviour with Artificial Neural Network and Multiple Linear Regression Models”, *Engineering Structures*, vol. 33, no. 3, pp. 903 – 910, 2011. DOI: 10.1016/j.engstruct.2010.12.011

J. Mata, A. Tavares de Castro, and J. S. da Costa, “Time–frequency Analysis for Concrete Dam Safety Control: Correlation between the Daily Variation of Structural Response and Air Temperature”, *Engineering Structures*, vol. 48, pp. 658–665, 2013.

K. P. Murphy, *Machine Learning: a Probabilistic Perspective*. The MIT Press, 2012.

L. H. Nguyen and J.-A. Goulet, “Anomaly Detection with the Switching Kalman Filter for Structural Health Monitoring”, *Structural Control and Health Monitoring*, pp. e2136–n/a, 2018. DOI: 10.1002/stc.2136

Y. Ni, X. Hua, K. Fan, and J. Ko, “Correlating Modal Properties with Temperature using Long-term Monitoring Data and Support Vector Machine Technique”, *Engineering Structures*, vol. 27, no. 12, pp. 1762–1773, 2005.

Y. Ni, X. Zhou, and J. Ko, “Experimental Investigation of Seismic Damage Identification using PCA-compressed Frequency Response Functions and Neural Networks”, *Journal of Sound and Vibration*, vol. 290, no. 1-2, pp. 242–263, 2006.

K. Pearson, “LIII. On Lines and Planes of Closest Fit to Systems of Points in Space”, *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 2, no. 11, pp. 559–572, 1901.

B. Peeters and G. De Roeck, “One-year Monitoring of the Z24-bridge: Environmental Effects versus Damage Events”, *Earthquake Engineering & Structural Dynamics*, vol. 30, no. 2, pp. 149–171, 2001.

- B. Peeters, J. Maeck, and G. De Roeck, “Dynamic Monitoring of the Z24-bridge: Separating Temperature Effects from Damage”, in *Proceedings of the European COST F3 Conference on System Identification and Structural Health Monitoring, Madrid, Spain*, 2000, pp. 377–386.
- G. E. Roberts, *From Music to Mathematics: Exploring the Connections*. JHU Press, 2016.
- S. Särkkä, *Bayesian Filtering and Smoothing*. Cambridge University Press, 2013, vol. 3.
- J. Shlens, “A Tutorial on Principal Component Analysis”, *CoRR*, vol. abs/1404.1100, 2014.
- M. Tatin, M. Briffaut, F. Dufour, A. Simon, and J.-P. Fabre, “Thermal Displacements of Concrete Dams: Accounting for Water Temperature in Statistical Models”, *Engineering Structures*, vol. 91, pp. 26 – 39, 2015. DOI: 10.1016/j.engstruct.2015.01.047
- G. Welch and G. Bishop, *An Introduction to the Kalman Filter*, Chapel Hill, USA, 2001.
- M. West and J. Harrison, *Bayesian Forecasting and Dynamic Models*, series Springer Series in Statistics. Springer New York, 1999.
- R. Westgate, K.-Y. Koo, and J. Brownjohn, “Effect of Solar Radiation on Suspension Bridge Performance”, *ASCE, Journal of Bridge Engineering*, vol. 20, no. 5, p. 04014077, 2014.
- K. Worden, G. Manson, and N. R. Fieller, “Damage Detection using Outlier Analysis”, *Journal of Sound and Vibration*, vol. 229, no. 3, pp. 647–667, 2000.
- Y. Xia, B. Chen, S. Weng, Y.-Q. Ni, and Y.-L. Xu, “Temperature Effect on Vibration Properties of Civil Structures: a Literature Review and Case Studies”, *Journal of Civil Structural Health Monitoring*, vol. 2, no. 1, pp. 29–46, 2012.
- Y. Xia, B. Chen, X.-q. Zhou, and Y.-l. Xu, “Damage Identification in Civil Engineering Structures utilizing PCA-compressed Residual Frequency Response Functions and Neural Network Ensembles”, *Structural Control and Health Monitoring*, vol. 18, no. 2, pp. 207–226, 2011.
- Y. Xia, Y.-L. Xu, Z.-L. Wei, H.-P. Zhu, and X.-Q. Zhou, “Variation of Structural Vibration Characteristics versus Non-uniform Temperature Distribution”, *Engineering Structures*, vol. 33, no. 1, pp. 146–153, 2011.
- A.-M. Yan, G. Kerschen, P. De Boe, and J.-C. Golinval, “Structural Damage Diagnosis under Varying Environmental Conditions—Part I: A Linear Analysis”, *Mechanical Systems and Signal Processing*, vol. 19, no. 4, pp. 847–864, 2005.

K.-V. Yuen and S.-C. Kuok, “Ambient Interference in Long-term Monitoring of Buildings”, *Engineering Structures*, vol. 32, no. 8, pp. 2379–2386, 2010.

K.-V. Yuen and S.-C. Kuok, “Modeling of Environmental Influence in Structural Health Assessment for Reinforced Concrete Buildings”, *Earthquake Engineering and Engineering Vibration*, vol. 9, no. 2, pp. 295–306, 2010.

APPENDIX A ENVIRONMENTAL EFFECT OBSERVATIONS MODEL FOR T3

Figure A.1 presents the hidden state variables time dependent marginal distribution for environmental effect observation $E_3 = T3$.

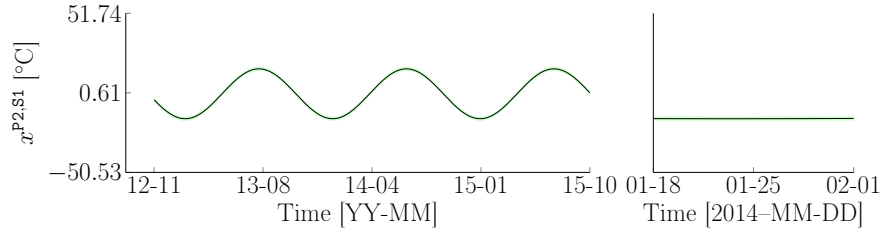
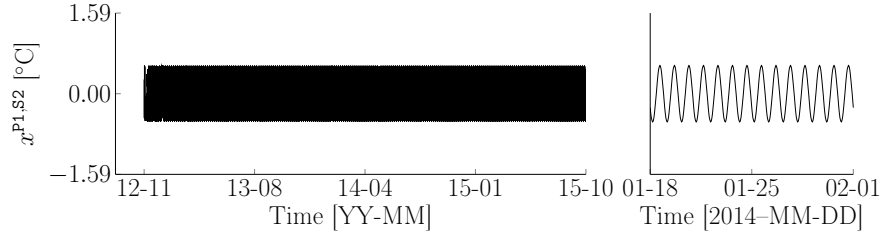
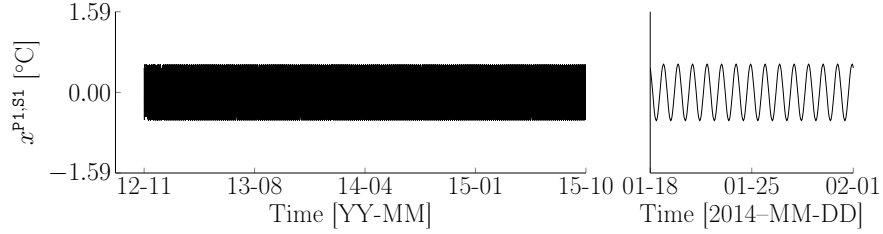
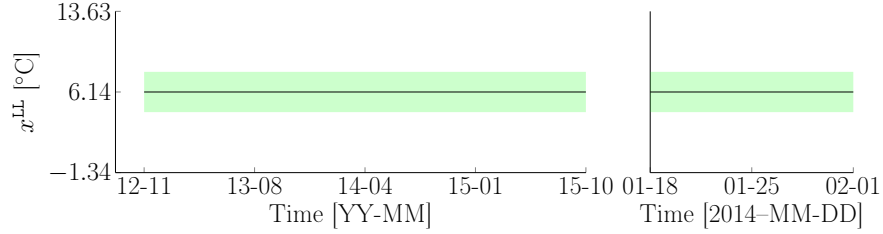
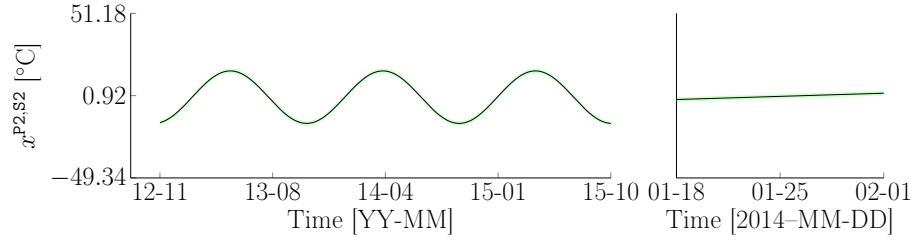
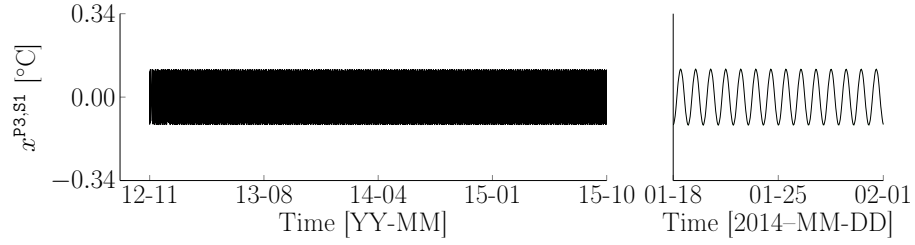


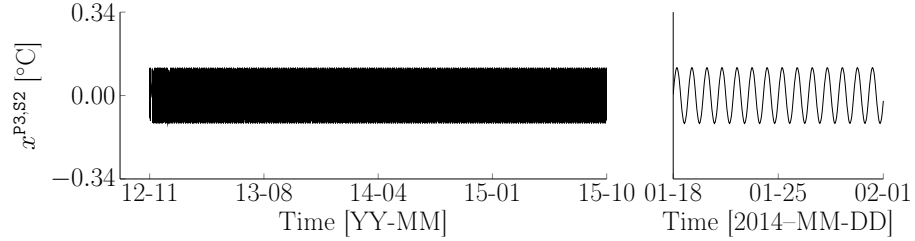
Figure A.1 Hidden state variables estimation for the model of temperature observation T3. The left and right parts show the distribution for the entire dataset and for a period of 14 days respectively



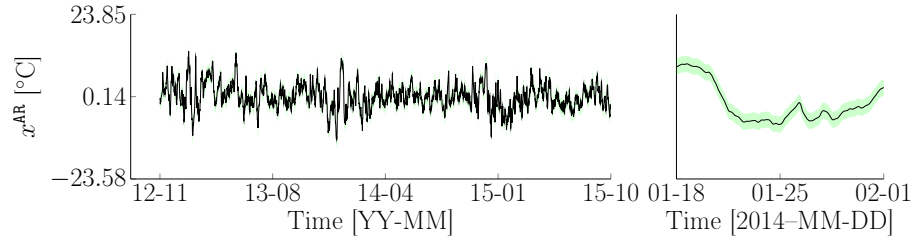
(e) Second hidden state variable of periodic component $P2 = 365.24$ days



(f) First hidden state variable of periodic component $P3 = 1.0027$ days



(g) Second hidden state variable of periodic component $P3 = 1.0027$ days



(h) Autoregressive component

Figure A.1 Hidden state variables estimation for the model of temperature observation T3. The left and right parts show the distribution for the entire dataset and for a period of 14 days respectively

Figure A.2 presents the the raw data and model prediction.

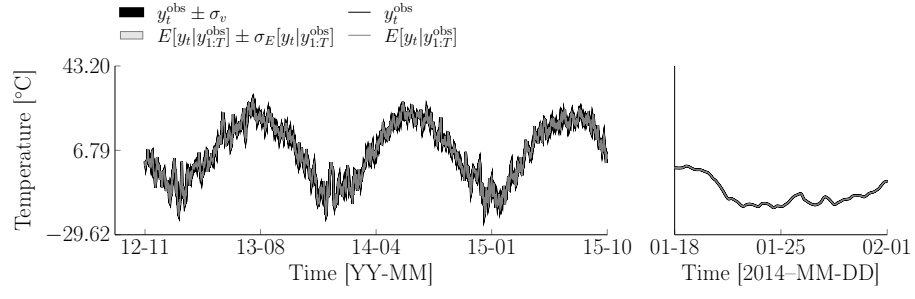


Figure A.2 Observed data (y_t^{obs}) and model prediction for temperature T3

APPENDIX B DISPLACEMENT MODEL FOR CASE 2

Figure B.1 presents the hidden state variables time dependent marginal distribution for structural response observation d of case 2.

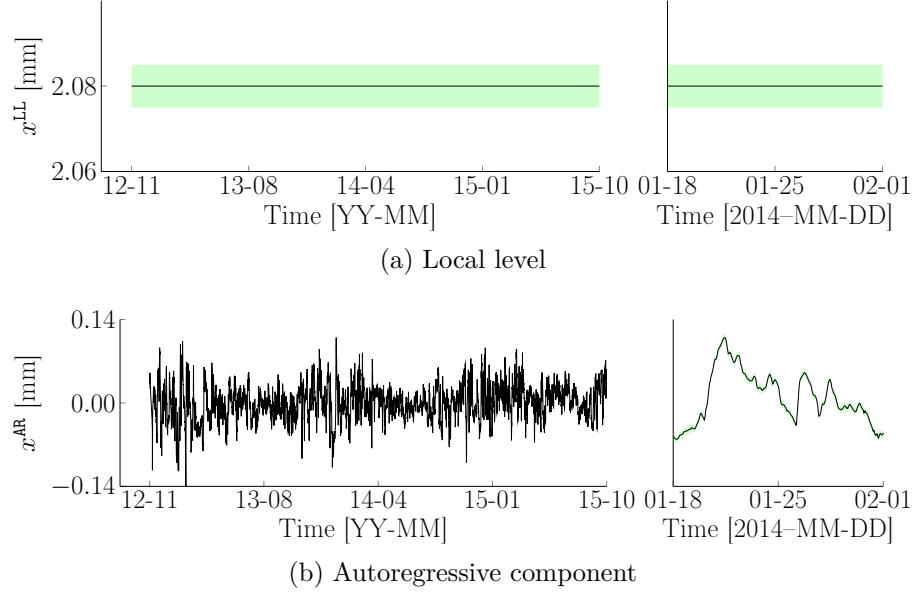


Figure B.1 Hidden state variables estimation for the displacement observation model of case 2

Figure B.2 presents the the raw data and model prediction.

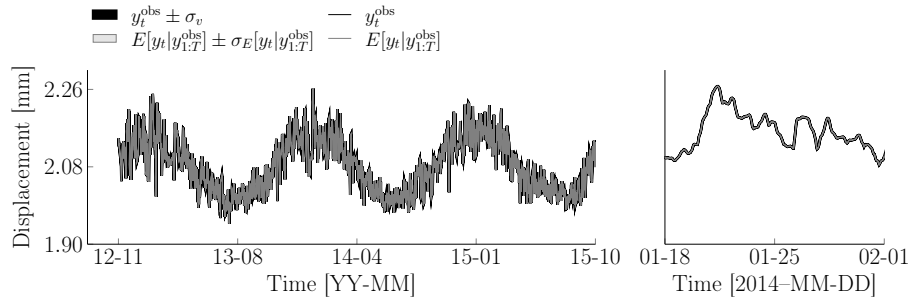


Figure B.2 Observed data (y_t^{obs}) and model prediction for displacement observation d from case 2

APPENDIX C DISPLACEMENT MODEL FOR CASE 8

Figure C.1 presents the hidden state variables time dependent marginal distribution for structural response observation **d** of case 8.

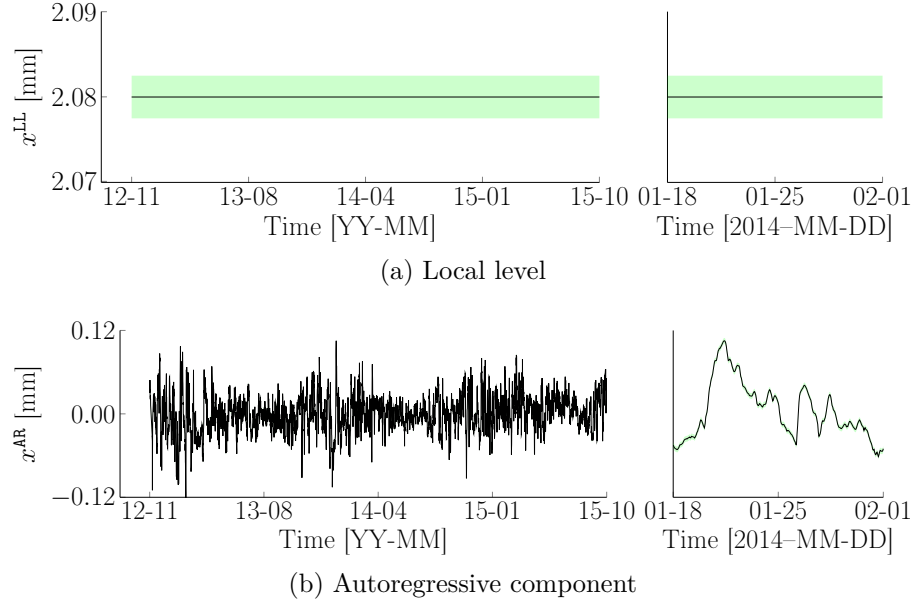


Figure C.1 Hidden state variables estimation for the displacement observation model of case 8

Figure C.2 presents the the raw data and model prediction.

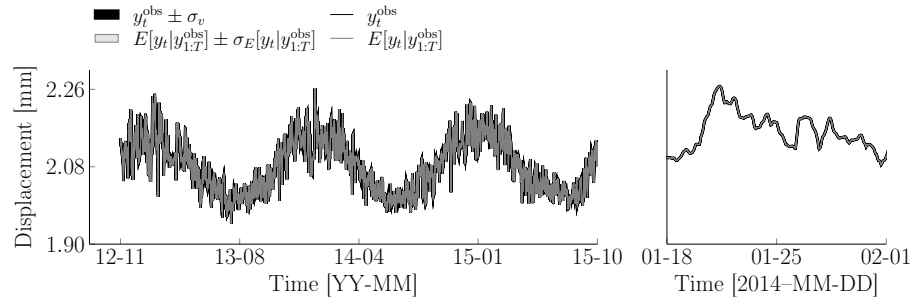


Figure C.2 Observed data (y_t^{obs}) and model prediction for displacement observation **d** from case 8